ON CONTROL OF MOTION PATHS AND SPEED IN A SPLINE-BASED
ANIMATION

BY

WOOSEOB JOU

# ACKNOWLEDGMENTS

ii

Finally, Thanks are due my wife Chongmi for the care of our home and our children through the years of the preparation of this dissertation.

# TABLE OF CONTENTS

iv

Abstract of Dissertation Presented to the Graduate School
of the University of Florida in Partial Fulfillment of the
Requirements for the Degree of Doctor of Philosophy

ON CONTROL OF MOTION PATH AND SPEED IN A SPLINE-BASED
ANIMATION

By

WOUISOR JOU

August, 1991

Chairman: Dr. John Staudhammer
Major Department: Electrical Engineering

Various graphical models for the visualization of pharmacologically-
active proteins have been developed to help recognize a three-dimensional
structure of molecules; these models have been principally aimed at
visualizing the same images of molecules. This work is an extension of such
computer representations in animated display of molecular activities such as
the folding of protein molecules. This dissertation describes a new graphical
model suitable for the animation of protein folding.

During animation, the motion path of individual atoms forming a
molecules may exhibit an unobservable wiggle or inflection under certain
conditions. The $C^2$ class of splines can be used to eliminate these artifacts, but
a mathematical representation of the $C^2$ splines has not been developed to

vi

data. We develop an analytic expression for the $G^1$ splines. Moreover, we extend and apply the $G^1$ splines for the generation of user descended camera path as an interactive animation environment.

This dissertation develops a method behind the incremental knot spacing to control the motion speed by an approximation in the discrete parameter domain so that control of motion speed is achieved. We acquire temporal aspect from spatial aspect of animation while providing a means of limiting key frames.

Based on these major tools for the control of motion path and speed, we develop a molecular animation model termed the *slot backbone* model, in which a molecular backbone is depicted as a combination of slanted cylindrical surfaces. More importantly, the interweaving technique was incorporated into the model to smooth the motion transition.

# CHAPTER 1
## INTRODUCTION

### Protein Folding

In late 1957, Christian B. Anfinsen, a researcher at the National Institutes of Health, made a remarkable discovery while experimenting with the only notion of protein folding that promotes the enzymes of 'enzymler' that more proteins to renature their native structure (Northl). What causes newly made proteins—which resemble loosely coiled springs—to wind into a specific shape?

The experiment confirmed that a more multitude of informations can fold onto a denatured structure and back into the exact original native structure (Beech). Under proper conditions, the denatured shape remembers its original native shape and completely recovers it. Therefore, it was generally believed that the extent and sequence of a protein, a one-dimensional trial, was fully sufficient to specify the molecule's ultimate three-dimensional shape and biological activity.

The next question to arise was "How does the protein fold between the two states?" For more than 30 years after Anfinsen's discovery, hundreds of investigators examined this problem (Chuloh, Wetlth) in an attempt to predict the folding pathway. The solution is of more than scientific concern. The major impetus behind the biotechnology industry is to produce new and useful proteins. It is already possible to design genes which could guide the synthesis of such proteins. For an essential feature of such synthesized

1

promise must be that they be able to retain the natural, stable state with full biological activity by subsequecy holding. If however, the experimental proteins do not fold or unreclad, the resulting proteins cannot exist in vivo.

Because only a role nestral is required for folding (e.g., 1-100 nm) and because the methods that provide structural information about protein folding (e.g., X-ray crystallography, Nuclear Magnetic Resonance) are notoriously slow (Kinell), so are but been able to exhaust the folding process. Instead, biochemists have managed to map a few folding intermediates in a stable form by slowing down the process. These folding intermediates constitute the basic building blocks for the various models which have been proposed to discuss the folding process (Kinell). However, there is as yet no known solution to the problem of prediction of the folding pathway (Kinell, Kirfl).

With the advent of the computer age, theoretical endeavors could complement earlier experimental work. For example, the shape of a folded protein ought in principle be determined by an empirical surface potential energy function. Given the structural values of the charges of atoms, the resulting energy may be expressed as the sum of pair interaction between all atoms of the protein (Kirfl, Kirfl). Factors that influence the protein-fold energy function (e.g., van der Waals forces, bond stretching energy, and chain torsional energy) depend heavily on the postulated model (Kirfl, Cirfl, Lurfl).

Based on these factors, the computer adjusts the coordinates of the atoms so that the overall energy is lowered until a minimum is reached. The minimum energy state is the natural, stable state of the protein. The process, known as energy minimization, has been a valuable tool not only for confirming models of structures but also for determining folding pathways by

mendations [Smit87, Chak91, Lev90b, Rei79]. Although some problems in this approach have been reported [Chak92], progress in this area of research has been rapid. This dissertation, however, does not deal with the modeling and computational processes involved in such energy stabilization.

Besides their use as a macromolecular task, protein sequences have found computers useful for the display of protein structures. Interactive computer graphics put bond structure, molecular volumes, and surface areas render the control of protein engineers with a simple interface devices [Pab94]. The protein images can be rotated, magnified and represented as even the user. Fragments of proteins can be manipulated to approach one another and joined to test the ease of fit of a substrate and a binding site. Given these the visualization of the binding interactions to the design of substrates that bind more effectively and enhance the activity of complex systems [Mit83]. In essence, computer graphics plays a key role in protein engineering.

Various graphical models [Lev84, Mun85, Rei79] for the representation of three-dimensional structures of proteins have been used to aid in the recognition of proteins—real or simulated.

The skeletal model is based on the Kendrew-Watson physical models of a molecule [Ric81] and shows the molecular framework which is a collection of lines joining atoms. The skeletal model, often referred to as wireframe model [litt94], implies atoms graphics by the joinlines of bonds and their center. Although the skeletal model has capability of rapid imaging of large molecules, recognition is restricted as intermolecular topology early to bond lengths, bond angles, and relative atomic positions.

The spherical model illustrates molecular surface shapes as formed by the union of van der Waals contact surfaces. Van der Waals surfaces, which are used as atomic radius in the spherical model, is defined as the distance

where a repulsive force begins to appear when two atoms are drawn together. Although the spherical model enables enhanced perception and conceptualisation of orbital interactions (Chem3D commonly occurring at the contact surface, existing graphics hardware does not permit interaction with a spherical model larger than a few tens of atoms [Ref9, Fu10].

The ribbon model (Cn46, Ref6) visualises the backbone of protein, ordered at the repetition at nitrogen, alpha-carbon, and carbonyl-carbon atoms. Common implementation of this model represents the backbone line as a ribbon composed of multiple parallel, smooth flattened turning along its length.

With the growing interest in protein folding, the obvious extension of such computer representation is an animated graphic display of the motion of a molecule. Feldmann and Levitt [Ref46], in a 1983 film of the molecular dynamics of bovine pancreatic trypsin inhibitor (BPTI), demonstrated the power of computer animation as a tool for the researcher. From a collection of moving pictures, all aspects of the reaction pathways could be seen directly. Particularised concepts were easy to grasp using visual effects involving colour changes. Furthermore, the application of computer animation shows promise as application ranging from user-interactive robotic surgery to the prediction of structure for physicochemically-active enzymes and the animated display of such molecular processes as DNA unwinding or protein folding process.

Nevertheless, these existing molecular models are appropriate for the representation of the static view of protein structure. In applying them for molecular animation, however, a number of serious problems arise. The flipping of side chains in the skeletal model distorts the viewer's attention. The spherical model, besides being prohibitively expensive even to produce images for individual folding intermediates, hides the backbone with the

outside confines of atomic spheres. This distinction deals with the definition of a new model which reveals those internal forces and is specifically designed for use in molecular animation.

## Statement of the Problem

### Considerations as a Body View

Energy minimization is a predominant simulation technique to pseudomotion a routine pathway, since large molecular movements which involve connected displacements of many atoms to those requiring the creating of energy barriers cannot be investigated with molecular dynamics alone [Ande, Levitt]. Because protein folding is associated with a plethora conformational change, molecular structures resulting from energy minimization have been widely used in modeling the folding, interactions [Cedd]. Our animation system also uses energy minimization as a tool for the prevention of folding intermediates.

A peptide plane in a protein is defined as the rigid plane containing carbon, oxygen, nitrogen and hydrogen atoms which exists between two adjacent alpha-carbons belonging to corresponding amino acids. Therefore, a protein molecule may be regarded as a sequence of the peptide planes joined by alpha-carbon atoms.

Because the bond lengths between atoms are fixed, only the angles joining each of the peptide planes can freely rotate. Consider the rotation of the peptide planes induced by a movement of a single alpha-carbon atom. A local rotation of a peptide plane propagates through the entire chain length to reduce the constraints imposed on the bond lengths. From a mechanical

viscogene, the situation is somewhat similar to the tension of the vein with respect to the body as a multi- none robot system 53d57. When only the joint between body and feet arm or closed to exists, while the rest of the arms have no degree of freedom to rotate, all of the arms will undergo nearly the same degree of rotation as the first arm.

Energy minimization, however, does not calculate the regular rotation of the peptide planes in a one-by-one manner. Rather, energy minimization adjusts the coordinates of a whole group of localized atoms in such a way that the overall energy is lowered and a minimum is reached, without concern for the individual rotation of the peptide planes. Thus when we see as a result of energy minimization in the location of atoms which have gone through conformations of rotators clustered at several joint positions.

The structures model has been used most frequently in previous attempts at molecular mechanics because of the high computation cost required when the spherical model is used. In this model, the solvent particles are represented by the junction of bonds and their torsion. Nevertheless, the structures model imitates the following fundamental forces at an imitation end.

As a result of the combination of rotators prolonged from several joint positions, the angles joining the bonds can vary randomly. What makes the animation possible is the virtual (feature of the free image restraining until the created image is displayed. If the free image continuously imitates from the previous one, and therefore continuous through the entire succession of frames, conservation of movement cannot occur. Since the angles joining the bonds do not cause gradual changes, the viewers cannot rely on data as a sense of visual character. The bonds belong the junctions under the same law of visual enjoyment as the worldwide model. Since the

junctions of bonds were is apparently random: way, the later channeling these junctions cannot show continuous movement. Figure 1-1 shows an animation sequence of part of a hemoglobin geodesic. Nowhere in this figure dissipate the finite sequence numbers. As can be noticed, it is hard to find a coherence object which is moving as a usually smooth motion, or this example of the nowhere-moded. Instead, each disturbing features on flipping side chains and abrupt changes in bond angles can take away the mover's attention during animation.

Figure 1-2 illustrates a ubiquitin molecule—whose abnormal structural deformation is known to be characteristic of Alzheimer's disease—shows in the spherical model. Besides the enormous computation involved in rendering the individual spheres, the backbone atoms are buried under an unusable contrast of spherical surfaces.

What remains local to correspondentin in the animation of fidelity is the global structural change represented by the backbone of protein. This can be evidenced by numerous illustrations [ske77, fle60, foo48] where the folding process is shown by thick and smooth-backbone lines. Defining backbones lines or the continuation of the later greatly the response of mimigens, alpha carbon and carbonyl-oxime atoms, backbone a priority to interpret the folding process in terms of the extremes and contractions of these backbone lines.

Characterization of the folding process is possible with the backbone lines. Because amino acids stick to each other to form proteins we call the individual amino acid the monomer, or primary structure. Newly-formed proteins are often coiled random state, implying that no region of backbone looks significantly different from other regions. Certain portions of these coils, which are probably unstable and fluctuating, may move to "stable" strand which stable coils will eventually coil. The backbones of proteins that

formed has some peculiarities. For example, the usual backbone of the nerve molecule can be divided into regions of secondary structure, which are distinct segments having a characteristic conformation. They may form before, far down, or turns interacting the before and stretch of sheet. The presence of different secondary structure makes the possibility that certain amino acids tend to be found in a specific secondary structure (RNA?) for instance, some amino acids are found more often in helices than in random coil backbone lines in the interior play an important role in a conceptualization of secondary structure.

An additional important advantage of using backbone lines to represent the molecule is that the backbone lines are completely recover geometrical relationships of major amino Protein's designers quite often require the ability to stop the minimum flow and query about the geometrical information because amino as a state sense of a some. Naturally they may not a amino to club as a portion of the backbone states which shown their structure during selection. Because the two peptide planes oriented in alpha-carbons are rigid units, docking are the portions of alpha-carbons can generate coordinates of other atoms on these rigid units. Provided with a proper data structure to link the alpha-carbons with the rigid units, these atoms can be viewed in the spherical model. Therefore, it is possible to view the globular changes in structure and to focus on and analyze the spherical shapes in a static sense. The linkage between the two modes is an most click on the location of the alpha-carbons.

The ribbon model has succeeded in helping the user visualize the backbone lines. This model, however, has some inherent problems. First, conventional implementations of the ribbon model do not pass through the backbone atoms. In Figure 1-3, the central ribbon of the ribbon is the primary

8

curve approximating the backbone line, and the remaining threads are constructed to be parallel to it. However, it is not assured that the control thread will fix the position of the stylus rather, a rough constraint of the backbone minus. Most of the visual implementations (Carib, Sphtf) estimated the smoothness property so highly that approximating stylus became distractant in reproducing the backbone line? In part, popularity of approximating stylus stems from the fact that the interpolating splines draw wriggles in markedly-curved helices. The approximating splines, to be described later in this chapter, are not guaranteed to pass through the sample points which in this case are the backbone curves. Since the threads of the ribbon model do not necessarily extend the location of the backbone curves, the position of the backbone cannot be marked on the backbone curve, and the curser loses the bridge between the estimation and point view

Second, in the ribbon model, it is difficult to recognize depth cues for the only clue to the three-dimensional understanding of backbone lines is the twist and width change at parallel lines of that ribbon model. Often, the threads are intertwining and the patterns as in which backbone conformation is clear: it left to the viewer's experience. The viewer needs time to sensitize such details to form a good understanding of the spatial nature of the backbone. In animation, however, clues are harder the direction of a static frame. The animation stream cannot arbitrarily assign a long time to figure out the three-dimensional structure depicted by multiple ribbon threads. Even if this model is used to make individual frames of the animation sequence, the multiple threads become another source of flipping, as was the knots in the wireframe model.

Figure 1-1. A schematic diagram of the annotation of a waveframe model
Black dots represent positions of an identical atom in different frames.

Figure 1.2 Spherical model of a silicprint molecule.



Figure 1.3 Ribbon model of the silicprint molecule in Figure 1.2. Wireframe model is overlapped in the background.

## Considerations in a Dynamic View

A general minimization of previous extraction systems to molecular graphics [Aab0, Fgb0, Tmj0, Sjx00] is that extraction may be performed by simply displaying the individual frames concurrently. If the individual frames are of a relatively similar spatial pattern, this assumption may be true. However, folding intermediates, whether they are determined from a theoretical production or they are taken from experimental results including a enzyme extraction such as energy minimization, need to be dissimilar. Therefore, the resulting frames exhibit abrupt and untoward displacement from one frame to the next. Such disruptions can be noticed in Figure 1-1 where the transition between Frame 1 and Frame 2 is conspicuous when compared with other transitions.

Key frames derive their naïve flaws not by drawings [Tmj00] as the pictorial measurements used as movement indication, all as the individual holding intermediates of the energy measurements in the key frame: we can employ the traditional computer graphic technique called inbetweening; inbetweening [Kor77, Kor81, Tle80] works to balance the estimator question two key drawings and the computer calculates and generates additional drawings between them by minimizing the distance between pairs of corresponding items. When applied to the region of folding intermediates, the technique will generate additional intergroup intermediates between the two holding intermediates. Movement from one holding intermediate to another can be improved by the introduction of additional intermediate frames. These inbetween frames identify the sharp conformational discontinuities between the folding intermediates three important.

what we are concerned with is the visualization of the global movements of a backbone structure, reionization of these inherent factors does not interfere as long as it aids in the expressing of motion.

Suppose an atom has to travel from point A to point B as a result of energy minimization such that point A belongs to the first binding intermediate and point B belongs to the next one. There can be various motion paths from A to B depending on the interconverting scheme. The path expectation might be straight lines or smooth curves (splines in a broad sense, can be said to be the smooth curves connecting points A and B. If the curves pass through the two points, the splines are said to be interpolating. In contrast if the curves approach and do not pass the two points, the splines are said to be approximating. For the smoothest it is bound, the class that fits the target destination & approximating splines are used. This means that the approximating splines used in interconverting code will modify the structure of the original binding intermediates, making active frames imaginary ones. If a user stops the animation scheme to analyze the static structure, he has no place to anchor because the original key frames do not belong to the animation sequence. Therefore, switching between animation and static analysis becomes difficult with the approximating interpolations.

In an attempt to pass through the key frames, interpolating splines may replace the approximating splines. The motion path passes through the exact instances of the original binding intermediates used as key drawings. The spline, however, introduces the problem of wiggles, loops, or undulation. Suppose that the atom proceeds to move from point B to another point C. If the distance from A to B is large compared with the distance between the points B and C as in Figure 4 or Figure 5, or is as in Figure 6-3, the

motion path tends to wiggle in the neutral [8, 9]. We would like to avoid this misbehavior of the motion path in our animation system.

In addition to the problems related to the motion path, there remains another problem associated with the control of motion speed. If an alien accelerates and decelerates unnaturarily transition between holding intermediates will become fascinating as well. Such motion speed is obtained from sampling the point types on the motion path. With more points sampled, more laborious frames are produced and the motion will be slower. The major factor determining the motion speed is the spatial distance between sample points. Traditionally, it has been known as a difficult problem to adjust the motion speed on a curve produced by a parametric representation. Because the spline used in computer graphics are predominantly expressed parametrically, the inner-frame distances should also be controlled by the parameter. However, the relation between the parameter and the inter-frame distance is not linear in general, and this nonlinearity complicates the control of the human speed.

## Objectives and Approach

We state our problems as follows:

Problem 1: The spherical model is not suitable for folding animation since the backbone conformation may be trouble.

Problem 2: The conventional mark is not suitable for the sensation because of its inherent discontinuities in spatial posture.

Problem 3: The ribbon model cannot be used directly for interaction because the viewer needs time to traverse the twist and width of the parallel lines

Problem 4: Since the details of the offline model do not pass through the backbone stores generally, information about those locations is lost.

Problem 5: Backhole curves shown with interpolating splines may show wiggles in some regions.

Problem 6: Sliding intermediates tend to be desirable.

Problem 7: The technique of lateravaraing by approximating splines can miss the key feature.

Problem 8: Inlaceraraing by interpolating splines can exhibit an abnormal motion path because of the wiggles.

Problem 9: Even if interpolating splines are used, motion speed may not be controlled.

This dissertation resolves Problems 1 through 3 by deriving a new model called the solid backbone model. The solid backbone model shares the basic idea of subtracting backbone lines with the ribbon model. The backbone can sufficiently conveys all the information about folding while reducing the time constraints at the spherical model. However, the solid backbone model explores the multiple threads of the ribbon model with a series of stacked cylinders to speed up the virtual integration process during animation.

Being defined as a team of cylinders with proper care on the backbone lines, the solid backbone model gives the feeling of three-dimensional depth by proper shading at the cylindrical surface: it is a known fact that, in animation, relatively rough shading will suffice, compared with static scenes. Since the novel of correct hardware is to facilitate rapid computation, generation of rough shading of three-dimensional objects at a speedy process: for instance, most current graphic devices are equipped with the VLSI circuits

dedicated to studying. Rendering features related to the solid backbone would not presented at Chapter 5.

Problem 4 can be handled by applying the refocusing techniques described. Such illuminating is a 'must' to enable the recognition of the smooth transition of folding pathways. Were it not for the role of refocusement themes, the entire motion sequence would be a discrete array of spatially condensed images. The refocusers feature help the viewer conceptualize the motion.

Problem 4 and problem 7, although apparently different, are practically the same in light of the interpolation scheme. If the approximating splines are applied to entire positions in space to produce backbone curves on a static scene, problem 4 arises. If the same approximating splines are applied to the motion part of a single atom to produce the path trajectories over a span of time, problem 7 arises. Since visual interactivity is extremely important in current applications of computer graphics, the transition between the tentalizer mode and the static mode should be ensured right from the beginning. Therefore this dissertation adopts interpolating splines both for the determination of static shapes of individual frames and for the determination of the motion path of individual atoms between the frames. That way, the backbone curves will pass through the backbone atoms and the motion path will let the key frames make up of folding intermediates.

The adoption of interpolating splines naturally invokes problem 5 and problem 6, which are coherent architectures of interpolating splines in space and time domains, respectively. Could we overcome the singular from the interpolating splines? Chapter 3 resolves these problems by developing a new class of interpolating splines termed free-knot splines.

Problem 4 is unavoidable whether we choose cryptosurfacing splines or interpolating splines. Could the motion speed be extended with care? Chapter 6 establishes an algorithm to control the motion speed on the path trajectories produced by spline techniques in general.

In summary, the state of the dissertation can be categorized as follows:

1. Development of negative interpolating splines.

2. Control of the motion speed in animation.

3. Development of a new model for molecular animation.

Chapter 5 deals with the first category while Chapter 4 handles the second. Chapter 3 describes the development of a prototype system as an implementation of the model as the final category.

## Organization

Chapter 2 explains fundamental concepts involved in the spline techniques used in computer graphics. Spline function in the context of parameter curves are stated. Visual continuity, one of the most important concepts used in developing our new class of splines, is described in that chapter.

Chapter 3 develops a new class of interpolating spline called free form splines, which are visually continuous. Chapter 5 is organized as follows:

1. Motivation in relation to molecular animation is presented.

2 A survey of the previous research associated with the control of wagpies is made.

3 Mathematical definition of general tree form spines is presented.

4 Mathematical definition of steadily continuous tree form spines is presented.

5 Four characteristic types of tree form spines are defined in terms of relevant tangent vectors.

6 Performances of the four types of spines is measured on the removal of the wagpies.

7 A comparison of the tree form spines and the conventional approach is made.

8 As an application of the tree form spines, procedures for an interactive design of motion path, are described.

Chryter 4 develops a method for controlling the machine speed in animation. It is organized as follows:

1 The motivation for control of the motion speed is described in terms of multimeder animation.

2 Research relevant to the control of speed is summarized.

3 As a framework of user control method, a method called dynamic hint spacing is defined.

4 A pair of enhancement methods are presented in a tool for further refinement of the dynamic time spacing.

$\delta$ The dynamic knot spacing and the enhancement methods on represented by a pseudo-code.

A finer analysis of the algorithms is presented.

Chapter 5 describes an application of the above techniques to the scripts of protein folding. This chapter develops a prototype system that uses the four have splices developed in Chapter 3, both for the rendering of the backbone and the motion path. The dynamic knot spacing developed in Chapter 4 is incorporated into the prototype system to control the motion speed in animation. Chapter 5 is organized as follows:

1 Following an introduction, general background knowledge of protein chemistry as applicable to the problem is summarized

2 A literature review of molecular graphics in general, and molecular animation in particular, is made

3 Fundamental concepts delineating our animation model are described

4 Considerations involved in the design of the prototype system are explained

5 Experimental feature resulting from the prototype system are discussed

Chapter 6 summarizes the major accomplishments and contributions of this dissertation

# CHAPTER 2
# PREREQUISITES AND REVIEW OF SPLINE FUNCTIONS

The replication of the act of parametric curves and surfaces can be viewed as the origin of Computer Aided Geometric Design (CAGD). Until the 1960's the only reliable tool for the communication between the designers and the manufacturers had been numbers. For instance, in a car-body engineering, the surface had to be tacitly expressed with numbers. Since the numbers were at best a discrete approximation of the continuous real shape, a set of curves was carved in a 3D model and actual interpretation to mold the shape was left to the experience of highly skilled patternmakers.

With the introduction of the bézier curve (Bézier, CAGD had a major breakthrough. Being geometrically represented, the bézier curves exhibit all data to be expressed exclusively by a few numbers. Since the specification of numbers defines the entire curve shape, the shape could be standardized regardless of the experience of the patternmaker. Taking advantage of the approach, mathematic became involved in generating a variety of curved schemes that could interpolate or approximate given numbers so that CAGD could be further applied to the design process in general.

This chapter describes the prerequisites terms and definition for such computations, and presents previous work done in this area of computer graphics to facilitate further discussion in subsequent chapters. Only those terms necessary for the exploration of the approach will be given and the coverage of previous work will be limited to this context. Research done in very specific areas will be addressed in later chapters.

22

Often, a purely mathematical definition of a term involving redundant symbols and terminologies makes a concept extremely difficult for graphic designers to understand. Therefore, preference is given to more concise yet precise representations of such a mathematical definition. For example, the matrix representation of a spline function is easier to understand [Smith] and thus has been more widely used than polynomial representation. The description in this chapter are given in such a way that complicated mathematical definitions are minimized.

## Parametric Curves

Although a curve may be generated by a collection of points, provided they are closely spaced, there are several reasons why a mathematical representation is popularly used in computer graphics. The advantage of a mathematical representation is that it is precise and the properties of curves such as slope can be easily calculated [later]. Moreover, it can be stored compactly as a computer in the form of equations.

Before a parametric or a nonparametric form can be used to represent a curve mathematically. An explicit nonparametric space curve is given by a set of equations of the form

$$x = x$$
$$y = f(x)$$
$$z = g(x)$$

where

$x$ and $y$ are arbitrary mappings from $x$ to $y$ and $z$, respectively. For instance, let $f$ be

$$ax^2 + bx + c$$

so that

$$y \cdot bx^2 + bx + c_i = d$$

where

a, b, c are terms arbitrarily coefficients to be solved

Given the coefficients of these data points on space, a is simple to get the a, b, and c value so that the x,y relationship of the curve passing through the points can be calculated. However, such nonparametric expressions are inadequate for the purpose of computer graphics. The reasons are as follows:

1. Frequently, the situation arises where coordinates of one data point and the slopes of two data points are specified instead of specifying three coordinates. If the slope is taken from a data point where the curve slope is perpendicular to the x-axis, the value of the slope will become infinity, making evaluation of the coefficients impossible.

2. If the curve shown above forms such a loops, there will be multiple values of y corresponding to a single x value, and the evaluation of the coefficients becomes difficult.

3. When points on a nonparametric curve are calculated with equal increments in x, the portion of the point will be distributed evenly along the length of the curve so that the quality and accuracy of the graphical output is affected.

These difficulties are obviated by auto-dependency, and can be remedied by using parametric representation commonly used in computer graphics. In parametric form, each coordinate of a point on a curve is represented as a function of one or more parameters. For a curve with one parameter, the position vector for a point on a curve is expressed by the parameter. For

notation, the three Cartesian components of a space curve in terms of a parameter $t$ are written as

$$x = f_1(t)$$
$$y = f_2(t)$$
$$z = f_3(t)$$

where

$f_1, f_2, f_3$ are arbitrary functions from $t$ to $x, y, z$, respectively.

Since a point on a parametric curve is specified by a single parameter value, the parametric form is scale-independent. For instance, the equation of a unit circle in its nonparametric form

$$x = 1$$
$$y = \pm\sqrt{1 - x^2}$$

can be converted into parametric form as

$$x = \cos\theta$$
$$y = \sin\theta$$

with

$$0 \le \theta \le 2\pi.$$

While single Cartesian basis uses to 3rd degrees, a closed circle can be generated without producing the dual $y$ values appearing in the nonparametric representation. The tangent vector at a point on the circle, with respect to the parameter $\theta$, is given by

$$P'(\theta) = (-\sin\theta, \cos\theta).$$

In contrast to the infinite tangent of the nonparametric form with no $y$ value of zero, the parametric tangent with a $\theta$ value of zero does exist and is not infinite.

These advantages of parametric representation have led to its popularity as a tool for computer graphics, as their most interpolation schemes

are developed with parametric notation. Our approximation regarding in the subsequent chapters will also be developed and explained with parametric representations.

## Spline Functions

The interpolation theory of mathematics has borrowed the term spline from its physical counterpart, a mechanical device used by draftsmen to fit a curve of minimum curvature through successive points of a set. Physically, a spline (Sheldt, Sayfti, Sexdl) is a flexible strip which is bent under the weight of heavy metals placed to control the resulting shape. By varying the number and position of the lead weights, the spline can be made to pass through the specified data points such that the resulting curve appears smooth or flat.

Given a mathematical point of view, the problem of defining a curve from a known set of data points is one of interpolation. For instance, a curve can be made to pass through all known points by use of polynomial interpolation. It is important to note that a clear distinction is made between the definition of the class (many)ulatory polynomial and the polynomial used as a spline. For instance, given ten data points, an interpolatory polynomial is a single polynomial passing through all ten points. On the other hand, a spline is a piecewise polynomial passing through a subset of the ten data points. The complete curve is produced by linking the adjacent polynomials at some point between the subsets. In the example, three of the splines, each passing through four consecutive data points, can be combined to form the final curve. In general, the mathematical spline is a piecewise polynomial of degree $k$ with continuity of derivatives of order $(k-1)$ at the

common point between segments. For instance, the cubic spline has second order continuity at the joints.

In practice, the function representing a spline does not have to be polynomial. It can be any function including, the exponential and trigonometric functions. However, the polynomial geared popularity as computer graphics because it is simple and easy to evaluate, and relatively differentiable on those points of a curve which do not belong to data points.

A curve generated by such pieces of splines, expressed in terms of the parametric representation, can be defined as a continuous map of a collection of intervals $[s_i, ..., s_{i+1}]$ into three-dimensional space Figure 3.1 illustrates the mapping from the knot variable $s$ onto the parameter $t$ and the mapping from $t$ onto the space curve $P(s)$. Initially, points $P_i$ through $P_n$ can be regarded as the points on the curve $P(s)$ corresponding to the parameter $t_1$ through $t_n$ to the curve Plot P this requires that the experimenter determine the three dimensional space. Each knot number $s_j$ in this sense is defined as a breakpoint or a knot and the collection of all $s_j$ is called the knot sequence Note that the lateral definition of the knots distinguishes knots from data points. The data points on the curve as three-dimensional vectors while the knots are the corresponding parameter value used for the generation of the space curve first interpolation then. A local parameter $u$ for the interval $[s_j, s_{j+1}]$ is defined by setting

$$u = \frac{s - s_j}{s_{j+1} - s_j}$$

so that the normalized parameter $u$ can range from zero to one while the global parameter $s$ varies from $s_0$ to $s_n$. It is by this mapping from the unnormalized parameter space used for the generation of the space curve. A local parameter then can be expressed in terms of the normalized parameter $0 \le u \le 1$. Because of

this mapping, one can unfortunately generate the regression for each piecewise polynomial without paying attention to the actual values of the knot parameter $u$. Subsequently, curve $P(u)$ is a collection of piecewise polynomials as a function of the parameter $u$.

Since the polynomial is infinitely differentiable on the curve, the problem of differentiability is concerned with the data points making up the joins between two adjacent piecewise splines. For instance, suppose a cubic spline $P_1(u)$ that interpolates two data points $P_0, P_3$ with normalized parameter $u$. Furthermore, assume that we know the derivatives at the points, $P_1$ and $P_2$. Then the equation of a cubic polynomial curve satisfying these four constraints is

$$P_i(u) = [u^3 \ u^2 \ u \ 1] \begin{bmatrix} 2 & -2 & 1 & 1 \\ -3 & 3 & -2 & -1 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} P_0 \\ P_3 \\ P_1 \\ P_2 \end{bmatrix} \qquad (31)$$

This equation is for one cubic spline segment and can be generalized for any two adjacent cubic segments. If additional cubic spline segment, $P_2(u)$ through $P_4(u)$ are assumed to interpolate internals $[D_3 P_4]$ through $[D_{n-1} P_n]$, the curve interpolating the points $P_i$ through $P_n$ is made up of $n$ subdivision of piecewise splines $P_i(u)$ through $P_n(u)$.

Let us apply the preceding definition of a mathematical cubic spline which is two times differentiable $C^2$ continuous at the joints. To request the second order continuity constraint at the joints, we are supposed to specify the positions of the $n$ data points and the two endpoint derivatives, $P_1$ and $P_n$. The first order derivatives or intermediate data points cannot be specified arbitrarily. Instead, they are calculated by the constraints of second order

combining the entries, if the number of data points r is equal to three, the derivative at the point $P_2$ is obtained by setting

$$\left.\frac{d^2 P_i(s)}{ds^2}\right|_{s=1} = 0$$

In general, the first order derivatives of the intermediate data points are computed by solving the tridiagonal matrix equation

$$
\begin{bmatrix}
4 & 1 & 0 & 0 \\
1 & 4 & 1 & 0 \\
0 & 1 & 4 & 1 \\
0 & 0 & 1 & 4
\end{bmatrix}
\begin{bmatrix}
\dot{P}_1 \\
\dot{P}_2 \\
\dot{P}_3 \\
\dot{P}_4
\end{bmatrix}
=
\begin{bmatrix}
3(P_2 - P_0) \\
3(P_3 - P_1) \\
3(P_4 - P_2) \\
3(P_5 - P_3)
\end{bmatrix}
\tag{2.3}
$$

where

$n$ is the the number of data points to be computed, and

$P_1$ and $P_n$ are the first and last data curve.

The coefficient matrix of this equation is diagonally dominant, the solution exists and is unique [Fa78]. The calculated first order derivatives are used to produce the corresponding piecewise splines in each interval.

Although the above mathematical cubic spline is twice differentiable at these points, it is not local and therefore, it is neither exact in computer graphics. That is, if a position of a point is modified, the entire system of linear equation should be evaluated again for all the first order derivatives. Consequently, the change of a single data point propagates and modifies entire piecewise splines. This loss of locality has kept the mathematical cubic spline from being widely used, since an attempt to modify only a portion of a curve will affect all other portions of a curve which are already completely

designed; in contrast, if the color option has only the contrary of first-order derived ms, the interpolating curve changes only in the vicinity of that point. The tradeoff for the lower differentiability in locality of the cytlass. Due research is mostly concerned with these bands of color splines which preserve locality.

The most common splines employed in computer graphics applications are local color splines, local because displacement of a knot affects curve shapes only in the vicinity of that point, color because third order or appropriate to get relatively smooth interpolations with reasonable computational complexity. Moreover, it is of the lowest degree spline curve which allows a point of influence and thus has the ability to twist through space [Bag76]. Local color splines can be categorized into interpolating splines which pass through their data points and approximating splines in which data points are used only to control the curve shape. To provide more smoothness near data points, the approximating splines have continuity of the second order derivative. The price paid is that the data points, with the exception of endpoints, are no longer on the curve. The data points act as the shape which they typically define, only remotely. Because of this property of the approximating splines, the data points are often called control points.

Because of the undesirable feature that the data points are not on the curve, approximating splines require practice in figuring out the resultant shapes when the sample points are changed. In general usage, the data points are imperfectly modeled on a trial and error basis until an acceptable shape is achieved. Because of this property, approximating splines are often used as a modeling tool for rough approximation of the curve and then the curve fit is done using interpolating splines.

Interpolating splines are referred to by several different names: the cardinal spline, the Catmull-Rom spline, or the Overhauser spline. Twentieth Century Fox's splines are often called cardinal splines. Interpolating splines pass through given data points and the property makes them ideal for our animation system despite their relative shortcoming that sample points have only the first order derivative continuity.

## Visual Continuity Conditions

This section describes a visually continuous class of splines which will be mathematically defined, proved, and extended in the next chapter.

Somewhere of a curve [Rog76, Lee86] is defined in terms of its differentiability from a mathematical point of view. A curve which is two times differentiable is smoother than one with the continuity at first order derivative over data points. As is shown with the continuity of first and second derivatives and that their magnitudes and directions are identical. To fulfill this requirement, individual derivatives at the $n_i$, $s$ components of Cartesian coordinates with respect to a parameter should be of the same magnitude.

Nevertheless, it is a known fact in computer graphics that the parametric tangents do not provide an appropriate measure of continuity in the sense of geometry [Far88, Bar88, Ros89]. For instance, Figure 2-2 shows two piecewise bits cubic splines given by

$$P_i(t) = [2t, 4t] \qquad 0 \le t \le 1$$

$$P_2(t) = \left[ 2 + \frac{1}{2} t, 4 + t \right] \qquad 0 \le t \le 1$$

and the parametric first derivative vectors are

$$P_1 \mid d = [2, 1]$$
$$P_2 \mid d = \left[\frac{1}{2}, \frac{1}{2}\right]$$

Even if the segments join with a discontinuous parametric first derivative vector, the joint at $[r, 3]$ is visually continuous. Figure 3-5 is another example illustrating the fact that the visual continuity does not necessarily coincide with the continuity of the parametric derivatives. The line segments are,

$$P_1(s) = [s(0s - 3^-, s(0s - 3^-)] \qquad (0 \le s \le 1)$$
$$P_2(s) = [s + 3(-1s)^2, 3(1 - 1^2)] \qquad (0 \le s \le 1)$$

Therefore, the parametric first derivative vectors become

$$P_1 \mid d = [s0 - 3s, -s0 - 3s] \qquad (0 \le s \le 1)$$
$$P_2 \mid d = [10 - s0s, -0s] \qquad (0 \le s \le 1)$$

Since both derivatives evaluated at the join are the same with the value of zero, the parametric first derivative is continuous at the joint positioned at $[s, 3]$. Mathematically, this continuity is called the $C^1$ continuity. Nevertheless, the two line segments at that point exhibit a cusp which is far from visually continuous.

Visual continuity is defined in terms of the continuity of the unit tangents. If the unit tangent of the left and right segment at a joint is identical, it is visually $G^1$ continuous. The unit tangents of the two segments at the join are given as

$$T_1(s) = P_1 \mid d \left[ \frac{-3}{\sqrt{3^2}}, \frac{-3}{\sqrt{3^2}} \right]$$

so that it is $G^1$ continuous. On the other hand, those of Figure 3-5 are

$$T_1(s) = \left[ \frac{s0 - 3s}{\sqrt{s^2 + s^2}}, \frac{-s0 - 3s}{\sqrt{s^2 + s^2}} \right] \qquad (0 \le s \le 1)$$

$$T_2(s) = \left[ \frac{10 - s0s}{\sqrt{1^2 + s^2}}, \frac{-0s}{\sqrt{1^2 + s^2}} \right] \qquad (0 \le s \le 1)$$

Except for the $z$ value of $2a$, it is not $G^2$ continuous by definition.

By defining visual continuity as the continuity of unit tangent vectors rather than the tangent vectors, the mathematical continuity somehow could conform to the visual perception. In addition to maintaining the same visual smoothness as the parametrically continuous splines, the class of the $G^2$ continuous splines encompasses a wider range of splines than $C^2$ continuous splines. These advantages of the visually continuous splines will be extensively exploited in Chapter 6.

## Literature Review

In approximation theory, mathematicians have developed various numerical analogs of splines, well before cubic splines were widely used in the computer graphics field. Among approximating functions, polynomial interpolation became a focus for study because polynomials are easy to evaluate, differentiate and integrate, and are well-behaved. As a result, classical (piecewise) polynomial techniques were innovated with finding a polynomial which, when interpolated between tabulated values of functions, could best recover the original function values.

The Lagrange form of interpolating polynomial $\ln(N, \ln(N))$ expressed in terms of parameters $z$ is given by

$$P(u) = \sum_{i=0}^{N} P_i \ L_i^N(u)$$

with the Lagrange polynomial

$$L_i^N(u) = \frac{\prod_{\substack{j=0 \\ j \neq i}}^{N} (u - z_j)}{\prod_{\substack{j=0 \\ j \neq i}}^{N} (z_i - z_j)}$$

where

$$\bar{P}_j \; (i = 0...s) \text{ are data points to be interpolated,}$$

$$\eta_j \; (i = 0...n) \text{ are knot values corresponding to } P_j \; (i = 0...n).$$

Despite its central formulation, Lagrange interpolation has few problems which limit its applicability in computer graphics. First, it oscillates in the order of the polynomial increases and the degree of this polynomial are not guaranteed to be preserved. This effect is often referred to as the Runge phenomenon [Far90]. Second, modification of one point affects the whole curve shape at the expense of costly computation. That is, the Lagrange polynomial lacks in locality. In most cases, this is not desirable for design purpose since modification of one portion of a curve must not jeopardize other completed portion by changing the entire shape.

The Bézier curve [Bez70, Bez74, Bez72] represents one of the earliest attempts to develop a flexible curve scheme in computer-aided design. Starting with de Casteljau's recursive algorithm [For87, For91], Bézier derived an equation for a piecewise approximating polynomial, which was used by Renault in design cor bodies. As adopted by Gordon [Gor74], the closed form of the Bézier curve is

$$R^n(u) = \sum_{i=0}^{n} P_i \, B_i^n(u)$$

with the Bernstein polynomial

$$B_i^n(u) = \binom{n}{i} u^i (1-u)^{n-i}$$

where

$n$ is the order of the Bézier curve,

$i$ is the normalized knot parameter,

$P_i \; (i = 0...n)$ are data points to be interpolated.

Notice that the curve passes through the two end points $\mathbf{P}_0$, $\mathbf{P}_3$ but not the external points. One of the important properties of the Bézier curve is the convex hull property [Far88]. Since the coefficients $B_i^n$ (u) are nonnegative, in addition to summing to one, they form a convex combination. A convex combination of points is always inside those points. This observation leads to the definition of the convex hull of a point set as the set that is formed by all convex combinations of a point set. More intuitively, the convex hull can be said to be a polygon bounded by the data points in the planar case.

A Bézier curve of the above form represents a smooth curve over a range of critical points. These geometric Bézier curves can be joined end-to-end to form a composite Bézier curve. At the joints of the individual curve segments, certain continuity conditions need to be met to preserve the overall smoothness of the composite curve. $C^0$ or higher continuity (CmN) in a composite Bézier curve requires that the change of control points satisfy certain constraints [Far88] at knot values. Therefore, if some critical point is modified arbitrarily to reshape curves adjacent to that point, the previous continuity condition of nearby curve segments may be destroyed.

The B-splines [Cor78b] were devised such that the same order of continuity relations after modification of control points. Geometric interpretation of B-splines in terms of the Bézier curve can be found in [Go88, Bart87]. Sharing the knot function with the Bézier curve, the B-splines inherit the convex-hull property [Bar93] from the Bézier curve Matrix representations [Bar93] of the cubic B-spline is

$$P(u) = \begin{bmatrix} u^3 & u^2 & u & 1 \end{bmatrix} \cdot \frac{1}{3} \begin{bmatrix} -1 & 3 & -3 & 1 \\ 3 & -6 & 3 & 0 \\ -3 & 0 & 3 & 0 \\ 1 & 4 & 1 & 0 \end{bmatrix} \begin{bmatrix} P_{i-1} \\ P_i \\ P_{i+1} \\ P_{i+2} \end{bmatrix}$$

where

$u$ is the parameter value at the range of $0 \le t \le 1$,

$P_{i-1}$ and $P_{i+2}$ are internal control points corresponding to $u = 0$ and $u = 1$,

$P_{i-1}$ and $P_{i+2}$ are internal points neighboring $P_i$ and $P_{i+1}$ respectively.

Noticing the continuity of parametric derivation is just a proper measure of visual continuity, Kestyú β spline (Barsh), Barsh incorporated a $G^1$
continuous (Barsh) tangent vector with left and second order continuity. It is important to note that the β spline belongs to the class of approximating spline. Hence the GR continuity is preserved at some points on the curve corresponding to the data points, but not to the data points themselves. Expressed explicitly, continuity of the two tangent vector and the curvature vector give

$$Q^{(1)}(0) = \beta_1 Q_i^{(1)}(u)$$
$$Q_i^{(2)}(0) = \beta_1^2 Q_i^{(2)}(u) + \beta_2 Q_i^{(1)}(u)$$

where

$Q_i^{(1)}(0)$ is the left order first derivative at vertex,

$Q_i^{(1)}(u)$ is the right order first derivative at vertex,

$Q_i^{(2)}(0)$ is the first order second derivative at vertex,

$Q_i^{(2)}(u)$ is the first order right derivative at vertex, and $β_1, β_2$ are shape control parameters

With the simple observation that a curve can preserve its visual continuity when the unit tangent vector rather than the tangent vector is identical for

the left and right dimensions, thereby could naturally deduce two parameters as a by-product. Upon inspection of these two parameters $\hat{y}_L, \hat{y}_R$ for their effect on the shape of a curve, they are referred to as two end-tension parameters, respectively. The line parameter controls the duration of the slope on a data point, while the tension parameters control the linearity of the curve segment, the more the parameter imposed the linearity of the curve.

The value $\delta$ affects and the $\beta$-spline are powerful modeling tools, they are able to model complex shapes easily. This "modeling" is carried out as an approximation process, manipulating the control points until a desired shape is achieved. However, the interpolating process should also be contemplated even if one sees approximating splines as a modeling tool. It takes the role of fixed refinement. The processes requirement that the curve pass through a given data set can even be imposed on the rough curve generated by the approximation process. While both β-spline and β-spline have disturbed complexity character to get the shape of a $C^2$ continuity basis, there have endeavors to deduce a class of interpolating splines which is capable of generating a smooth interpolating function with $C^2$ continuity only.

Catmull [CaP3] called the "Catmull-Rom spline" by separating blending function [Blm70] from the unified functions being blended. This concept applies to any blending function and to any method function to produce a wide class of splines. As one example of this class of splines, blending linear cardinal function with a β-spline base function yields a cubic interpolating polynomial. This polynomial, as source mentions, becomes

$$\mathbf{P}(t) = \begin{bmatrix} t^3 & t^2 & t & 1 \end{bmatrix} \begin{bmatrix} -1 & 3 & -3 & 1 \\ 2 & -5 & 4 & -1 \\ -1 & 0 & 1 & 0 \\ 0 & 2 & 0 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{P}_{i-1} \\ \mathbf{P}_i \\ \mathbf{P}_{i+1} \\ \mathbf{P}_{i+2} \end{bmatrix}$$

where

$\mathbf{P}_k, k = 0, 1 \dots 4 = OR$ are the control points,

$t$ is the parameter value in the range $[0 \le t \le 1]$.

As was remarked, this cubic polynomial maintains compulsatory character at the expense of $C^2$ continuity.

The above special class of Catmull-Rom spline was not first to appear. Brewer and Anderson [Bre77] surposed and described a lesser known spline called the Overhauser curve [Ove68] which was developed well before the Catmull-Rom spline. The Overhauser method, developed at Ford Motor Company in 1965, resorted to parametric blending function without any explicit knowledge of splining or of B-splines. Given four consecutive data points, $\mathbf{P}_1, \mathbf{P}_2, \mathbf{P}_3, \mathbf{P}_4$ position it can be made to manipulate the points $\mathbf{P}_2, \mathbf{P}_3, \mathbf{P}_4$. Similarly another position $\mathbf{P}$ can interpolate $\mathbf{P}_2, \mathbf{P}_3, \mathbf{P}_4$. Hence the linear blending of the parabolic A and B yields the cubic spline which is effective on the region $[\mathbf{P}_2, \mathbf{P}_3]$. This method forms the platform for the development of the new class of splines discussed in Chapter 3.

Olsen designers need to specify the slopes instead of interpolation points. In the case of the cubic spline, for example, they may enforce positions and slopes at two internal points instead of positions at four points. The question is, "Given data points and corresponding tangent vectors, find a $C^1$ piecewise cubic polynomial that interpolates the given data."

Cubic Hermite interpolation [Jac87, Par90], when applied to the end

plate derivative coefficient of Bézier curves, takes the following Bézier form.

$$P(t) = B_1 P_0^3(u) + P_{t=1} P_1^3(u)$$
$$+ (P_1 - \frac{D_1}{3}) P_2^3(u) + (P_{t=1} - \frac{D_2}{3}) P_3^3(u)$$

where

$P_1$ and $P_{t=1}$ are two end points of the region of concern,

$D_1$ and $D_{t=1}$ are slopes corresponding to $P_0$ and $P_1$,

$B_i^3(u)$ $i = 0...3$ is the Bernstein polynomial on parameter $u$.

In matrix form, the above equation simplifies to Equation 2.1 with $P_i$ replaced

by $B_i$. Interestingly, this formulation of the Hermite polynomial includes the

Central-Bias spline as a proper set. Central-Bias spline fabs into the special

case of Hermite polynomial if

$$D_1 = (P_{t=1} - P_0) / 2 \text{ and}$$
$$D_{t=1} = (P_{t=1} - P_0) / 2.$$

Kochanek [Koe84] noted this similarity and inserted multiplicative constants

to the expression for $D_1$ and linked this the resulting changes in the curve

shapes. Variation in these multiplications constants has led to define such

parameters as bias, tension and continuity. This method is compared in detail

with our strategy for the control of curve shapes in Chapter 3.

Figure 2-1. Mapping from a pericenter into a pericenter, and from a pericenter into the pericenter affine curve $P_0I$. The normalized parameter $s$ varies between zero and one while the curve runs from $P_0$ to $P_1$.
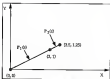
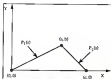Figura 2-2 Decovoluzione parametrica distruttiva per visually continuous case



Figura 2-3 Continuous parametrica decovolve per visually discontinuous case

## CHAPTER 3
## THE FREE FORM STUDIES

### Introduction

The development of computer-aided generative design relies on a wide range of mathematical methods for curve and surface fitting. The need for fitting curves and surfaces arises principally from the fact that many physical phenomena are continuous, although our measurement of them is discrete. From the discrete information, we try to reconstruct the continuum using the mathematical tool called spline interpolation.

In our situation, the position of points is represented by numbers constitutes each discrete information. Consider an atom and of a large molecule. During the folding process, the atom will change its position in a function of time in itshelves itself as dictated by the folding dynamics. However, because of the difficulties involved in the measurement, only a limited number of intermediate scenes can be captured experimentally. Therefore, the distance between the intermediate scenes will become unreasonably long.

The same holds true if such a folding process is simulated by a computer technique called energy minimization. For each step of the computation, the method tries to find a more stable atomic position by minimizing the sum of inter-atomic potentials, and the output of each step makes up a specific folding intermediate. As long as the new position yields less potential, the technique will reposition the given atom. Although the intermediate scenes are in more continuous if the distance is kept short, no

45

specific manner to given about the distance an atom travels for each step of computation.

There are two aspects in the evaluation of protein folding that require interpolation. First, various paths of individual atoms which make up folding intermediates should be interpolated. If an atom at position A at time $t_1$ moves into position B at time $t_2$, interpolation is needed to produce more animation frames of the atomic positions between the time interval (A,B). Second, in a strict sense in a fixed time, interpolation can be used to generate a backbone curve which passes through the positions of the backbone atoms. Because our animation system aims at viewing the folding process in the secondary structure level represented by a backbone curve, it is imperative to interpolate a given sequence of atoms. Therefore, an interpolation scheme could be used both for the generation of a motion path of an atom and the generation of a backbone curve in a static scene.

For reasons to be explained in Chapter 5, our animation model is based on the two principles in view of the interpolation scheme. First, the curve representing the motion path should pass through the folding intermediates. In terms of the above example, the generated curve should be both positions A and B. Second, the curve representing the backbone should also pass through all the backbone atoms. Therefore, usage of the interpolating spline rather than the approximating spline is essential for our animation system.

The ordinary cubic cardinal spline, a traditional interpolating spline which has been widely used, has an inherent undesirable *local property* known as a wiggle. Figure 3-1 shows the motion path of a backbone atom generated by a cardinal spline. In the blue curve[?] $b_{c1}(t)$, the motion path abnormally deviates from a reasonable track. Often, this type of curve behavior is called the wiggle. Due to the popularity of cardinal splines, there

as movement expressions describing similar behaviors such as no movement expression, or oscillation, as a resultant, a loci, or a structure of the curve. In addition to this orbital descent, the motion path of the retinol [$b_Z$] or merely flavor, thus making the movement of the main optic unaltered. While the curve shows a wiggle to a relatively short interval, it tends to be linear in the larger interval. The linearity illustration one of the mutual advantages of value optima, they drop our surface reflection and twist so that the curve looks more smooth and natural. Although retinol optima do not necessarily exhibit this property at all times, the inherent tendency to wiggle in short intervals tend to linearize at the larger interval is always present as amplified as the distance between the materials makes unduly. The same phenomenon can happen in the prevention of the backbone curves. If the first atomic positions in Figure 4-1 are considered in four different waves in the space domain on a given time, the backbone curves representing these atoms will suffer a similar problem, showing a loop and a line segment. This problem must be resolved for our system to work properly since the atomic positions during the folding process are not known.

This chapter is mostly concerned with the removal of the wiggles appearing in the ordinary retinol optima. An example of one approach is illustrated in Figure 3-2. The wiggle during the time interval [$t_0$ $t_1$] is removed by an interpolation scheme called the type-1 low interpolation. Moreover, the linearity at the interval [$t_4$ $t_5$] can be alleviated by this scheme. More importantly, the tendency to wiggle and linearize becomes unnoticeable, since it is based on a class of splines with $C^2$ continuity rather than the $C^1$ splines, to which the ordinary retinol splines belong.

position of the
stylus ckson
in three-dimensional space

a weight

an stylus ckson

motion path generated by an
ordinary cardinal spline
method

$t_0$   $t_1$   $t_2$   $t_3$   $t_4$   time $t$

Figure 5-1. Interpolation by an ordinary cardinal spline. A curve representing the motion path of an stylus ckson is shown with a thick line. The stylus ckson is a constituent of the backbone ckson. Each $t_i$ indicates the signal time during the folding process.
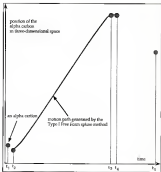
44



Figure 3-5 Interpolation by Type I spline showing the removal of the wiggle and aberration at the knotting. This class of spline ensures the G² continuity at interior points.

## Literature Review

Wiggins can be regarded as an expression of intense interest in the sphere. If the wiggins are properly controlled, they will lead to a visually pleasing reflection and twist. For instance, if the tendency to wiggle at Figure 2-1 is distributed to the longer interval, and the tendency to intersect is routed to the shorter interval, there remains no problem with the mirror scheme. The control of tension has been investigated in this arena.

An approach to relieve the inclination to wiggle by the adjustment of interval can be found numerous times in the literature. The most prevalent one, the tensioned cardinal spline (Smith), incorporates a tension parameter into the expression for the ordinary cardinal spline. The tensioned cardinal spline can straighten out the wiggles, or the excessive oscillation, by applying high tension to that region. However, while the wiggles in the shorter interval are straightened out by this method, the curve in a relatively longer interval tends to be looser or stiff. The fundamental flaw of this method is that the tension value cannot be applied for each piecewise interval because of the $C^1$ constraint.

Schoenberg's cubic spline (Schä, (S)IV, S) is called the spline under tension. Being $C^2$ continuous, the curve enables piecewise application of the tension parameter so that at each interval, the relative closeness can be consistently be reflected into the corresponding tension parameter. The only objection to this method is that it employs exponential functions instead of polynomials, thereby making evaluation of the interpolant more expensive and, in some cases, even unstable. Nielson's v-spline (Har9, (Sar9), Sar9) is a polynomial alternative to the Schoenberg's exponential function. The

orthochromatic, the method lacks the mathematical proof of the exactness of a unique solution [Burrel, Pratt] to the latter system of equations that leads to the value spline.

Foley [Foller, Sohn] developed a $C^2$ continuous cubic spline with the capability of piecewise tension control. Given the retained weights at each interval, the method finds a cubic polynomial that maximizes the sum of weighted measure of curvature, subject to the continuity conditions. Hence if a weight is set to be high in an interval, the curve will be assigned a low curvature so that it will be linearized. Since each interval can be assigned a different parameter a piecewise control of the curve is possible. However, a change of weight in one interval propagates through the entire interval to this method. The linear system of equations subject to the minimization constraint should be reevaluated even with a change of the position of a single data point. Therefore, the method lacks local controllability.

The $C^2$ continuity encompasses a wider range of curve shapes from the $C^0$ data, and contains the $C^0$ class as a proper subset. The principal reason why the $C^2$ class splines contains more diversified curve shapes is that the magnitude of the left and right tangents at a point is allowed to be different. Usually, the ratio of the left tangent and the right tangent takes the role of tension at that class of splines. Since the magnitude of the tangents directly affects the tension [Ward], the tension of the curve segments to the left and to the right of the point can independently be controlled in the $C^2$ splines.

In a sense, previous endeavors to control the tension while maintaining $C^2$ continuity at joints are unduly strained. A consequence of the $C^2$ continuity can be replaced by the right definition of the term tension. If softer represents making a curve appear to human eyes data points to follows data points. Subject to a $C^2$ constraint, the curve maintained next data

points (NerA, XaoB) cannot ensure the tension between two points. By the same token, the curve mentioned between data points (XaoA, XaoB, XaoC) cannot guarantee the similar degree of tension near the data points in contrast, as the $G^1$ sphere, the tension of the curve segments near the data points and between the data points can be controlled independently by allowing the magnitude of the left and right tangents to vary.

The properties of the $G^1$ continuity, the tonal continuity, is investigated in [MacD] and the application of this type of continuity in the approximating spline appears in [HarG], [RanW]. The only endeavor to produce a $G^1$ minimal spline is in [TiMSb]. The fact that the the basis functions of the minimal spline decay generally, a Catmull-Rom spline [CatM] can be expressed in terms of a Lagrange polynomials [ScA]. LaoM1 as opposed to the segments. Subsequently, the future minimal versions are termed parametrically (through an algorithm) for the piecewise Lagrange curves. Although the existence of such central motion is given without proof in its approach, by using such algorithms to structure subintervals [HarG], LaoM1 or in Catrom-s algorithm [Far58], these central version data lead to the future curve, which in this case corresponds to the $G^1$ central spline. Nevertheless, this approach is purely algorithmic "owing to the algebraic complexity" [TiMSb] involved.

The explicit closed form of the $G^1$ method cubic nodal can be produced by this approach. Generally, a disadvantage of such algorithm-provided approach is that a sufficient knowledge on the algorithm is required for its implementation. As long as a graphical programmer may understand it, it belongs to a complicated algorithm and it remains away from common usage. Moreover, the most fundamental advantage of using the $G^1$ spline: the control of tension by separate adjustment of the left and right tangents,

cannot be fully exploited in this approach because no specific elements was given to the importance of largest texture.

The SF autoregulates rather explained in this chapter will be derived to explain closed form starting with simple linear blocking of a pair of stress. Since the spline presented in this chapter is represented by a single unique matrix, evaluation of the curve becomes a simple matter. Furthermore, the tension parameters at this approach are closely related to the tangents at data points, so they provoke tension values can be intuitively adjusted.

## Mathematical Derivation of the General Nux Stave Spline

Local autoregulating spline have several characteristics which make them useful for the design process. Because the curve follows and passes through the control points, it is possible to inlink directly where to put a new point. Despite this constraints, piecewise sum of the verified splines have been limited because of the lack of proper tension control mechanism.

The familiarmind flow that closes the weights lies in the derivation of the cardinal splines. According to our observations, which will be shown shortly, the exact form direct correspondence on the slope of a curve. Between uniformisation moderately bit relatively considerable at inverting equal knot spacing condition, parametric splines which functional over at above form mathematician inversely followed this protocol. For instance, Overhauser [Pama69, Dou69] assumed that the knot value of a midpoint is the average at two adjacent knots. A similar assumption appears in Catmull-Rom spline [Cat74] by letting unevenry on the cardinal formula.

A rabbit confined within a linear blending of two quadrants. Each quadratic can be represented as the linear blending of two consecutive line segments where each line segment, in turn, is a linear blending of two points.

Let us start with the linear blending of two points. Given points $P_1$ and $P_2$, let us denote the linear blending of the two points $B$ simply the line segment joining $P_1$ and $P_2$. Points on this line can be represented in a parametrized form

$$f(u) = (1-u)P_1 + uP_2$$

where

$u$ is a parameter in the range $0 \le u \le 1$.

$f(u)$ is the positional vector of a point on the line.

For instance, point $P_1$ is the value of the function $f$ evaluated at zero while point $P_2$ is the value of the function $f$ evaluated at one. Therefore, as the parameter $u$ runs from zero to one, the value of the function $f$ makes up the continuous line segment from $P_1$ to $P_2$.

In Figure 3-8, $f(u)$ and $g (u)$ are two parametric functions which identify consecutive points $P_1, P_2,$ and $P_2, P_3$ respectively. The time parameter $u$ runs from zero to one making up the line segment $P_1P_2$, while the time parameter $t$ runs from zero to one making up the line segment $P_2P_3$. Therefore,

$$P_1 = f(0)$$
$$P_2 = f(1) = g(0)$$
$$P_3 = g(1)$$

Note that the two time parameters are independent of each other.

We wish to blend the two line segments linearly so that resulting curve becomes a quadratic parabola passing through the three data points. Let us denote the parametric representations of this parabola as $m(u)$. In order to

mentality (FaiNi) the paramters $e_i$ are let $F_x$, for the functions are evaluated at zero and let $F_y$ be the function as evaluated at one

$$F_1 = \frac{\partial f}{\partial y}$$

$$F_2 = \frac{\partial f}{\partial x}$$

The quantum is, "what is the least value corresponding to midpoint $P_2$?" Since point $P_2$ is located somewhere on the curve joining $P_1$ and $P_3$, the corresponding least value should be between zero and one. However, the exact least value is not known, and there is a degree of freedom in the assignment of the least value corresponding to the midpoint $P_2$, because the quadrature at hand will eventually be likened to produce a cubic spline, and because the three slopes of the cubic spline heavily depends on the functions being blended, the least value of the midpoint must be determined.

Before going further, we shall digress briefly to investigate the way infinitely related splines are derived. After an observation is made regarding the relation between the least value and the resulting curve shape, our own derivation will be resumed. This way, the concepts involved in the derivation of the free form splines can be clarified with ease.

Ordinary cardinal splines regard $P_2$ as a point mapped from the center of the parameter a distance such that

$a = 1$ and $a = \frac{1}{2}$

$b = ($ such $a = \frac{1}{2})$

In order for this assumption to be valid, the linear relationship among these parameters should be

$a = 2b$

$b = ja\ l$

Hence, at the arbitrary rational sphere, the quadratic parabola produced by the linear blending of the two line segments is

$$n(t) = (1-t)^2 l_1(t) + t^2 l_2(t)$$

$$= (1-t)^2 [(1-t) r_0 P_0 + t r_1 P_1] + t^2 [(1-t) r_1 P_1 + t r_2 P_2]$$

$$= (1-t)^2 [(1-t) r_0 P_0 + t r_1 P_1] + t^2 [(1-t) r_1 P_1 + t r_2 P_2]$$

$$= [a t^2, t, 1] \begin{bmatrix} 2 & -2 & 0 \\ -3 & 4 & 0 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} P_0 \\ P_1 \\ P_2 \end{bmatrix}$$

Let us examine the change in the parabola by varying the knot value. Our observation can be explained by Figure 3-6. Figure 3-6 shows the convergence of assigning two different knot values to the midpoint $P_1$ during the formulation of the blending. If the knot value corresponding to point $P_1$ is set to $R$, as is the case with ordinary rational splines, the parameter $t$ element is divided to the ratio of one to one. The resulting curve $m_1(t)$ exhibits a high inflexion in the interval $[P_0 P_1]$. However, with the division ratio of 1 : 4, the corresponding curve $m_2(t)$ shows an explicit trend to concentration such differences at the degree of inflexion. The curve shapes depend heavily upon the midpoint knot values, and we may offer the following.

I: Suppose we are bending a wire with two hands. The wire will bend more sharply if the distance between the two hands gets closer, provided that the wire lies fixed length. An extreme case of this bending is the origin of the wire. On the contrary, as the two hands are moving apart, the shape of the wire gets closer to a straight line. The hands working force will need to be relieved and tend to be loose. The linearity of a curve is proportional to the distance between output points.

I Consider two drivers assigned four and two hours respectively for delivering a package from Crossanville to Jacksonville. The driver assigned four hours must spend more time on the road than the one who is assigned two hours. That, the times at the former drivers will tend to encroach, while the times at the latter driver tends to be fewer. In the extreme case, the former driver may have a long between the two sides. If Crossanville and Jacksonville are two focal points to be triangulated, the travel set at the time constraints assigned to travel the distance between the two sides. That is, the intensity of a curve is inversely proportional to the assigned local values during the blending.

In order to allow the role at the locus as time constraints, we can make the local values proportional to the distance between the sample points. That is we allocate more time to travel in a greater distance.

In this section, however, we incorporate the local values of midpoints as parameters of the quadratics and will delay the determination of the values until a value spline is practised by linearly blending the quadratics. To do so other certified spline that we let the parameters values change to look into their initial positions, by allowing the local values to vary, our approach can be explicitly convinced with the arbitrary control points.

Let us examine the derivation of our free form splines by referring to Figure 3.0. To state the free value constituted as a new parameter, we introduce a parameter $u$ that indicates the local value corresponding to the midpoint $P_i$ (we will call this parameter the midpoint local parameter or midpoint local behaviour, since point $P_i$ is always in the middle of the two midpoints in terms of the sequence of data points). Therefore, point $P_i$ is the value of function as evaluated at a node that the parameter $u$ splits the parameter $s$ by the division rate at $u$ $(0 \leq u)$. Hence,

still = $F_1$,    mid = $F_2$,    half = $F_3$

That is, we think of point $P_2$ as the point mapped from a value of the parameter s passan, so that

$$s = 0 \text{ at } s = 0$$
$$s = 1 \text{ at } s = u$$
$$s = \tfrac{1}{4} \text{ at } s = \tfrac{1}{2}$$
$$s = \tfrac{1}{4} \text{ at } s = 1$$

Assuming a linear relationship among three parameters, the solution can be represented as

$$s = k_1 \, rsk_2$$
$$s = k_1 \, rsk_2$$

The constants $k_i$ of these equations can be determined by applying the above condition as follows that

$$k_1 = \frac{1}{u}$$
$$k_2 = \frac{1}{2(1-u)} = \frac{m}{1-u}$$

Therefore, the quadratic curves interpreting the new parameter u at the midpoint knot value u,

$$P_2(t) = G \, sl(t) w_3 \, (t)$$
$$= G \, sl(t) sP_{12} sM_{3} srw_3 (t) sP_{12} srw_{3} (t)$$

$$= [sl \ s \ 1] \begin{bmatrix} \dfrac{1}{u} & \dfrac{1}{u} & \dfrac{1}{1-u} \\ -\dfrac{1}{u} & \dfrac{1}{u} & \dfrac{1}{1-u} \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} P_1 \\ P_2 \\ P_3 \end{bmatrix}$$    (3.1)

The quadratic curve of the ordinary cardioid spline belongs to the special case of the class of quadratics where the midpoint knot parameter is between $\frac{1}{2}$.

Notice that the above quadratic has two different parameters: the parameter $t$ is the original parameter which moves between zero and one to make up a curve, while the parameter $s$ represents the location of the midpoint knot when leading to the variation of the entire curve shape.

We can view as a pointers to blend two of these quadratics to form a cubic spline.

As at Figure 3-9, consider two parametric quadratics mild and wild, each defined in terms of parameters $s$ and $t$ respectively, where $s$ and $t$ are in the range of 0 to a 1. Assume that mild is a quadratic which interpolates three points $P_1, P_2$ and $P_3$, such that the parameter $s$ varies between zero and one while the quadratic wild makes up the curve from $P_2$ to $P_4$. Similarly, the parameter $t$ is then used to interpolate three points $P_2, P_3$ and $P_4$. The parameter is now from zero to one respectively while the quadratic wild shapes the curve from $P_2$ to $P_4$. Also assume that the quadratic curves mild and wild are formalized with the midpoint knot parameters $s$ and $t$ respectively, so that point $P_3$ is the value at formula as evaluated at $s$, and point $P_3$ is the value of formula as evaluated at $t$. Therefore, the parameter $s$ splits the parameter $s$ by a division ratio of $s : (1-s)$, while the parameter $t$ splits the parameter $t$ by a division ratio of $t : (1-t)$. Hence,

$$mild = P_1 \quad mild = P_2 \quad mild = P_3$$
$$wild = P_2 \quad wild = P_3 \quad wild = P_4$$

Our aim is to derive an expression for the interpolating curve $pod$ which is defined in the interval $[P_2, P_3]$. We involve the individual quadratics symbols,

$$mild = (at + b) (B, P_1, P_2, P_3)^T$$

such that

$$B = \begin{bmatrix} \frac{1}{2} & \frac{1}{\sqrt{2}} & \frac{1}{2} \\ -i\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & i\frac{1}{\sqrt{2}} \\ \frac{1}{2} & -\frac{1}{\sqrt{2}} & \frac{1}{2} \end{bmatrix}$$

and

$$\text{old} = \{b\}^t \, U \, C \, \{P_x, \, P_y, \, P_z\}^t$$

with matrix

$$C = \begin{bmatrix} \frac{1}{2} & \frac{1}{\sqrt{2}} & \frac{1}{2} \\ -i\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & i\frac{1}{\sqrt{2}} \\ \frac{1}{2} & -\frac{1}{\sqrt{2}} & \frac{1}{2} \end{bmatrix}$$

If p(x) is a linear blending of $m(x)$ and $ed(x)$, it can be represented as

$$p(x) = (1-t) \, m(x) + t \, ed(x)$$

The linear relationship among parameters a, b and c is expressed as

$$a = c_1 \, ob_1$$

$$b = c_2 \, ob_2$$

where the constants $c_i$ are yet to be determined. The constraint between the parameters get

$$a = u + c = d$$

$$a = t + u = b$$

$$d = u + c = d$$

$$b = v + u = t$$

Applying some constraints to the above equations, we get

$$a = (1+d) \, (+u, \, and)$$

$h = v_0$.

Consequently,

$$p(u) = \Omega \cdot u \,\omega(du \cdot u \,\omega(u)$$

$$= (1 \cdot u)\, u^2 + \tfrac{2}{3}\, u\,(u+1) \in \|C$$

$$= (1 \cdot u)\,(1 \cdot u)d^2 \cdots \in D \text{ three } \|(4u + l|u)l \in \|C$$

Substitution of the matrices B and C into this equation yields cubic spline grid in terms of the midpoint knot parameters, $u$ and $v$. Hence we can define a new class of cardinal splines as follows:

*Definition:* The free knots of cardinal splines (abbreviated as the free form uniform bsosplines) are defined analytically as

$$p(u) = (u^3 \; u^2 \; u \; 1)\, A \,(P_1 \; P_2 \; P_3)^T$$

with

$$A = \begin{pmatrix} -v + \frac{1}{u} & 1 + \frac{1}{u} v + & v - \frac{1}{u} & v - 1 - \frac{1}{u} \\ 2v - 1 + \frac{2}{u} & 2v - \frac{3}{u} & -2v + 1 + \frac{2}{u} & 2v - \frac{1}{u} \\ -v & 0 & v & 0 \\ \frac{1}{u} & 1 - \frac{2}{u} & \frac{1}{u} & 0 \end{pmatrix} \tag{3.2}$$

where

$P_1$, $P_2$ $P_3$ and $P_4$ are four consecutive data points to be interpolated,

$u$, $v$ are the midpoint knot parameters corresponding to point $P_2$ and $P_3$ respectively,

$u$ is the parameter varying between zero and one, while the curve $p(u)$ ranges from point $P_2$ to point $P_3$.

The free item options include the ordinary combined options as a subset in which the parameter values are, $\nu = \nu' = \frac{1}{2}$.

Notice that the curve shape in the interval $[P_1, P_2]$ is determined by a set of $n$, $\nu$ parameters while that of the interval $[P_2, P_3]$ is determined by another set of $n$, $\nu$ parameters. The parameter $P_2$ may indicate the true value of the endpoint $P_2$ during the evaluation of the curve within $[P_1, P_2]$. Suppose we are evaluating the curves in the next interval $[P_2, P_3]$ by introducing another point $P_2$ following $P_2$. That time, the level value of the endpoint $P_2$ should be represented as $\nu$, and the level value of the endpoint $P_2$ takes the role of $\nu'$.

Therefore, the each successive evaluation of preemptive options, endpoint level parameters associated with each data point take the role of $\nu$ and $\nu$, in turn.

Instead of saying that the free item options are a direct blending of quadratics, one can regard the free item options as a cubic blending of data points. The cubic blending function, in this sense, is the multiplication of the matrix [a? of a ?] and a Campanile/xy, the elements of the matrix [a? of a ?] multiplied by the matrix A compose the four basic functions (Bezier, Go?/Bs) of the free item option. Notice that the sum of each row of matrix A is zero except for the fourth row, which is one. This means the normalization constraint $G(x?)$ required for the basic function.

It is most important to note that our derivation is strictly analytic in that the spline can be represented by a single four by four matrix. It is this simplicity of representation that underlies easy implementation of the spline functions in various design applications.

Insummary, we started out with the linear blending of two points. During the blending of the two consecutive linear segments, an extra degree of freedom is found in accepting the knot values corresponding to the

malignant. Based on the observation that the location of the knot values affects the shape of the quadratic parabola, we incorporated the endpoint knot values as new parameters into the formulation of the quadratic parabola. These parabolas, in turn, are blended linearly to yield the cubic spline defined as the free knot spline.

## Mathematical Description of the G2 Free Turn Splines

In the previous section, we created a new class of interpolating splines called free form splines. In an extensive design environment, the curve shapes which can be produced by the free form splines are almost unlimited. This variety mainly comes from the fact that the left segment and the right tangents are defined and controlled separately so that the shapes can be discontinuous at data points. However, we do not want to destroy shape continuity at data points at our objective. Since what we use for a curve, which are wiggle-free within intervals between given points, the splines interpolating backbones at the motion grab must remain visually smooth at the data points.

In this section, we extend the free form splines to an important class of splines called G2 splines, to create visual smoothness at data points. This is possible by differentiating the equations for free form splines and by imposing some constraints on the tangents that prevented from the observation of the formation of the segment of the ordinary cardinal splines, at a special case of free form splines, the concept of visual continuity is incorporated into the free form splines. A numerous candidate is chosen and verified for the free form splines to be visually continuous.

The free form splines belong to one of the widest class of cardinal splines. The term 'wide,' in this sense, means the variety of the curve shapes which can be generated by the splines. By varying the midpoint knot parameters (i.e. $a$ and $\varepsilon$), we can supple numerous curve shapes interpolating given data points. From the analyator's point of view, each variety leads to diversity in planning the median path between data points. On the other hand, return or derequers may use the variety to generate a multitude of curve shapes passing through the given data points in a sonic scene. However, the variety is useful only if we can control it.

The principal difference between free form splines and ordinary cardinal splines can be conceived in terms of the parametric derivatives at the data points. At each data point, the left and right tangents of ordinary cardinal splines always agree in both direction and magnitude. In contrast, free form splines distinguish the left tangent from the right tangent, and the direction and magnitude of the left tangent may differ from those of the right tangent.

The tangential properties of the free form splines can be described by fundamental notation. Suppose we are concerned in the estimation of the tangents at point $P_i$ (see Figure 3.5 for example). Let us denote the curves to the interval $[P_i, P_j]$ we need to inspect knot parameters $\gamma_i$ and $\tau_j$ corresponding to the point $P_i$ and $P_j$, respectively. Taking the parametric derivative of the free form splines, we obtain:

$$\frac{dz(u)}{du} = 3u^2 \Delta_i + 2u_1 \left( P_1 \ P_2 \ P_3 \ P_4 \right)$$

with

$$A = \begin{bmatrix} -a_0 + b_0^2 & -c_2^1 + c_1^2 & a_2 + b_2^1 & -c_2 + b_1^2 c_2^1 \\ 2a_1 + b_0^1 & -b_1^2 + \ldots & 2a_2 + b_2^1 + \ldots & b_1^1 + \ldots \\ -a_2 + b_0^1 & -b_2^1 & a_1 & 0 \\ b_2^1 & \ldots & 1 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} (2.3)$$

Hence the left tangent of point $P_0$ is given by

$$\frac{d\mathbf{r}(t)}{dt}\Big|_{t=0} = \left[ \begin{array}{ccc} a_2 + 1 & 2 \frac{1}{2} c_2^1 & c_2 \frac{1}{2}(c_1^1 + \ldots) \end{array} \right] [P_1 \; P_2 \; P_3]^T$$

By introducing another point $P_4$ same to point $P_4$, we can derive the right tangent at point $P_3$. If the two subjacent base parameters corresponding to $P_3$ and $P_4$ are denoted as $\tau_3$ and $\tau_4$ respectively, the right tangent of point $P_3$ is expressed as,

$$\frac{d\mathbf{r}(t)}{dt}\Big|_{t=1} = \left[ \begin{array}{ccc} -\psi_3(\tau_3) & \ldots & \psi_3(\tau_4) \end{array} \right] [P_1 \; P_2 \; P_3 \; P_4]^T$$

where

$\psi_i(\tau)$ represents the splint in the interval $[t_3,t_4]$.

Let us describe the ordinary cardinal splines on a special case of the free form spline. The ordinary cardinal splines replace both midpoint base parameters $\tau_3$ and $\tau_4$ with $\frac{1}{2}$, so that the left and right tangents agree. That is,

$$\frac{d\mathbf{r}(t)}{dt}\Big|_{t=1} = \frac{d\mathbf{r}(t)}{dt}\Big|_{t=0} = \frac{1}{2} [P_3 - P_1]$$

Geometrically speaking, the left and right tangents of the ordinary cardinal spline have the same magnitude and direction at point $P_2$. It is this correspondence of the left and right tangents as the ordinary cardinal spline

that is said to be $C^0$ continuous. More specifically, the tangent at point $P_0$ is one half the vector $(P_1 - P_{-1})$. Therefore, the magnitude and direction of the tangential vectors are fixed, once the data points are given; there is no way to control the magnitude and direction of the tangential vectors. As a result, the arbitrary cardinal splines are often estimated to be unduly attained. [think?]

In contrast, the tangent of the free form spline can vary depending on the selected endpoint knot parameters, as can be seen by the above equations for the tangent. As a matter of fact, one can assign two different knot parameters $v_0$ and $v_1$ for a single endpoint $P_0$, and produce curves with slope discontinuity at that point. In other words, the free form splines exclude those curves which have slope discontinuity $C^0$ intrinsic [?]

Although the free form splines are developed such that they would incorporate flexible tangents at data points, the same flexibility might hamper the smoothness. Splines having two different tangential directions at a data point cannot be said to be smooth. How do we impose the smoothness criterion on the free form splines while avoiding the unduly attained, reinvented statement of the ordinary cardinal splines?

Geometric continuity [Joe88, Far88] addresses the problem of parametric continuity. Parametric continuity means that the left and right derivatives with respect to the knot parameter are identical. However, as was mentioned in Chapter 2, the parametric continuity does not necessarily agree with our visual continuity. This is so because the parametric continuity is determined by the parametric derivatives are identical means that the derivatives of $P$ with respect to $v$, that is $dP/dv$, match at the common point. Nevertheless, the visual perception is based on the continuity of the tangential directions which are represented by a relationship between each tangent. Therefore, the visual

continuity, also referred to as geometric continuity, is preserved if the direction of the unit tangent is identical at a data point.

In light of the degree of the versatility in the curve shapes, $C^0$ of the free form splines covers the entire class of the splines including discontinuous curves. The $G^1$ splines lie in the middle of $C^0$ and $C^1$ in terms of the degree of the versatility. They are not so restricted as the $C^1$ splines, since they relax the tangent restraints of the ordinary cardinal splines. They allow the left and right tangents to vary in magnitude as long as both tangents have the same mid-tangent vector. Moreover, the direction of the tangents of the $G^1$ splines are not bound so the the vector from the previous points to the next point.

By definition [Bart87, Barsk88], the geometric first-order derivative is defined to be continuous iff continuous if

$$\text{Unit tangent at } \left.\frac{dQ(t)}{dt}\right|_{t=1} = \alpha \cdot \text{ Unit tangent at } \left[\frac{dQ(t)}{dt}\right]_{t=0}$$

where

$t$ is an arbitrary multiplication constant

Figure 3-6 shows two different tangent knot parameters $v_k$ and $v_k$ assigned to point $P_k$. Hence the evaluation of the curve in the interval $[P_k, P_{k+1}]$ takes the role of midpoint knot parameter so that the parabola A will be blended first-only, the parameter $v_k$ will make the parabola B be blended during the evaluation of the curve in the interval $[P_k, P_{k-1}]$. As parabola every point can be associated with two different midpoint knot values. The net effect of the dual definitions of the midpoint knot values is the ability to change the direction of the tangents at the data point. For example in the final video free

lions splines, $Dv_1$ is the sweeping tangent and $Dv_2$ is the sweeping tangent at point $P_2$.

Noticing that the $G^1$ splines overlook visual continuity while needing the tangent constraints of ordinary cardinal splines, we now impose certain constraints on the free form splines to make them $G^2$ continuous.

### Claim

Define four arbitrary endpoint knot parameters $u_1$, $v_1$, $u_2$, $v_2$ corresponding to points $P_0$, $P_1$, $P_2$ and $P_3$ as in Figure 54, such that point $P_1$ has two midpoint knot values associated with it. Then the free form spline are $G^2$ continuous if

$$u_1 = u_2$$

Thus the midpoint knot parameters are consistent for the evaluation of both the left and right splines as a sufficient condition for the free form splines to be $G^2$ continuous. In terms of Figure 54, the free form splines will be $G^2$ continuous if either the quadratic products A or B is used consistently, to evaluate the cubic spline in both the intervals $[P_0P_2]$ and $[P_1P_3]$.

**Proof**

Without loss of generality, we may treat as definition the first point as discussed earlier. The left tangent of point $P_1$ is,

$$\frac{dP(t)}{dt}\bigg|_{u=1} = \frac{1}{2}v_1(1-v_1)^2 P_0 + (1-2v_1+\frac{3}{2}v_1^2)P_1$$
$$+ (v_1-2v_1^2+\frac{3}{2}v_1^3)P_2 + (v_1^2-\frac{1}{2}v_1^3)P_3$$
$$= \frac{1}{2}v_1(1-v_1)^2 P_0 + a_1 P_1 + b_1 P_2 + c_1 P_3$$

where $a_1$, $b_1$, $c_1$ are the coefficients of $P_1$, $P_2$ and $P_3$ respectively.

$$= \frac{v_1^2}{2} [\mathbf{P}_2 \mathbf{P}_0] + v_1 [\mathbf{P}_1 \mathbf{P}_0]$$  (2.4)

The major tangent at the same point is,

$$
\begin{aligned}
\frac{d\mathbf{P}(t)}{dt}\bigg|_{t=0} &= (v_2 u_2 \frac{1}{2}) \mathbf{P}_1 + (v_2 u_2 \frac{1}{2}) \mathbf{P}_2 v_2 u_2 \mathbf{P}_2 \\
&= u_2 [\mathbf{P}_1 \mathbf{P}_2] + v_2 [\mathbf{P}_2 \mathbf{P}_0] \\
&= u_2 ([\mathbf{P}_1 \mathbf{P}_0] + [\mathbf{P}_0 \mathbf{P}_2]) + v_2 [\mathbf{P}_2 \mathbf{P}_0] \\
&= u_2 [\mathbf{P}_1 \mathbf{P}_0] - \frac{(2 u_2 v_2^2)}{2} [\mathbf{P}_2 \mathbf{P}_0]
\end{aligned}
$$  (2.5)

If

$$v_1 = u_2 = k$$

then

$$\frac{d\mathbf{P}(t)}{dt}\bigg|_{t=0} = \frac{d\mathbf{P}(t)}{dt}\bigg|_{t=0} = (\frac{k v_2}{2})$$

Since the right tangent is a scalar multiple of the left tangent by the factor of $(\frac{k v_2}{2})$ both tangents have the same unit tangent vector. This completes the proof.

Figure 5-3 Linear blending of two line segments. The line segments f(s) and g(s) are one linear blending of two data points $P_1$, $P_2$ and $P_3$ respectively. A quadratic in the interval $[P_1, P_2]$ can be produced by varying the normalized parameter s from zero to one inclusively.

Figures 3-4. Variation of the knot values corresponding to the midpoint $P_2$. The curve $w_0$ and $w_5$ are generated by inserting the midpoint knot parameter value of .35 and .5 respectively



Figures 3-5. Linear labeling of two quadrants: the curve represented by $P10$ is the cubic spline resulting from the blending of the quadratun mid and infit. It is defined in the region $(P_1, P_2)$

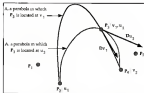Figure 3-4: The tangents created by two different midpoint knot parameters defined at a single point $P_b$. $Dv_1$ represents the tangent at point $P_b$ when its midpoint knot parameter is consistently equal to $v_1$, throughout the interpolation of the triples $[P_a, P_b]$ and $[P_b, P_c]$. Similarly the $Dv_2$ is the midpoint knot parameter has dual value at this point, the left and right tangents do not match as shown

## Vector Analysis of the C7 Rear Sutra Splines

In the previous section, we derived the relation between the left and right tangents of the C7 free form spline. The free form splines originally derived in studies require this interpolation but in common associative multiport knot parameter values, to assure valied smoothness at data points. Now we proceed to designate the exact values of the knot parameters to finally enable anglefree interpolation. If it is an interactive design environment, the parameters, and hence the curve shapes, can be varied over a satisfactory way in visual. However, one prefers to automatically deriving the besthave lines or the motion path empirically in interruptions by the designer's association. We do not want the flexibility to enforce the bestlieve lines iteratively until the curve shape is satisfactory for our association system.

This section analyzes the C7 free form splines by an iterative usage of vector calculus. Based upon the characteristic behavior of individual curves, three types of C7 splines will be developed. However, the application of C7 splines is not limited to these conception. Almost limitless kinds of C7 curve shapes can be produced in one approach, by varying a characteristic constant $\beta$, which will be explored shortly.

## C7 Continuity

Before delving into the C7 free form splines, let us digress briefly into the C7 free form splines to emphasize that a curve can be controlled by tuning the component vectors.

In general, the lless these splices can be said to be CP consistent in terms of continuity resolution at the data points. The left and right margins are not at the same locations at a data point, or can be evidenced by Equation 3.4 and 3.5. Therefore, in the general first order splines, the continuity condition that regulates the ordinary cardinal splines has been removed.

One of the advantages of repeater control of the left and right margins is that at being able to make an error-based step. For instance, the direction of a bouncing ball [Seckl] shows discontinuity near the instant point, and sometimes an estimator needs to simulate a means guide for such movements. More specifically, the motion path at the correct point is continuous but the tangents in that point are discontinuous. This sudden directional change can really be incorporated into the first form splines with proper control of the co-joined knot parameters. With deliberate assignment of these parameters, we could control even the level of the discontinuity.

Figure 3.7 illustrates the discontinuity effect. The curve marked $C1$, $C0$, and $C2$ represent three different instances of continuity variation, caused by assigning dual multiple knot parameter values to page $P_4$. The symmetrical tangents at panel $P_4$ can be the synthesized derivatives used in the last section. From Equation 3.4 and 3.5, we have

left tangent = $\dfrac{dp[u]}{du}\Big|_{u=1} = \dfrac{T_4^L}{2}(P_5-P_3)$

right tangent = $\dfrac{dp[u]}{du}\Big|_{u=0} = \dfrac{T_4^R}{2}(P_5-P_3)$

If we let

$T_4^L = T_4^R$

the $(P_4 P_5)$ component of the left tangent is the same as the $(P_3 P_4)$ component of the right tangent. Similarly, the $(P_3 P_4)$ component of the left tangent is

identical to the $(\mathbf{P_R}\cdot\mathbf{P_T})$ component of the right tangent. Since the magnitudes of the two vectors are identical, and since the magnitude of components vectors are symmetrical each other, the resulting tangent vectors are symmetrical with regard to point $\mathbf{P_T}$.

In practice, there can be in-between control strategies that could produce various shapes, and this example shows that the behavior of the component vectors can be a useful tool for the design and control of various path configuring the free form splines.

### Type 10/5 Free Form Splines

Let us consider the expression for the tangents of the G7 free form splines with 4-forward notations, such that

$$
\begin{aligned}
\mathbf{D_L} &= \frac{3}{4}(\mathbf{P_R}\,\mathbf{P_T}) + 0.14\,(\mathbf{P_T}\,\mathbf{P_T}) \\
\mathbf{D_R} &= b\,(\mathbf{P_R}\,\mathbf{P_T}) + \frac{(1-b)^2}{4}(\mathbf{P_T}\,\mathbf{P_T})
\end{aligned}
\tag{3.4}
$$

where

$\mathbf{D_L}$ is the left tangent at point $\mathbf{P_T}$.

$\mathbf{D_R}$ is the right tangent at point $\mathbf{P_R}$.

$b$ is the local parameter corresponding to midpoint $\mathbf{P_T}$.

Notice that the G7 free form splines have only a single local parameter associated with each data point, while the free form splines have two local parameters for each data point.

As was seen in the last section, the magnitude of the right tangent at a nodal endpoint of the left tangent by the factor of $\frac{(1-b)^2}{4}$. For example, if the midpoint local parameter is decreases, the multiplication factor increases, and the magnitude of the right tangent becomes greater. Physically, the

consequence of the decrease in tangential magnitude means that the curve
tends to persist in a certain direction for longer intervals. On the other hand,
if the midpoint knot parameter $t$ increases, the curve to the right of the data
point quickly loses the inclination to persist in the previous direction and
tends to form an inflection toward the next data point as soon as it leaves the
data point. Therefore, it is possible to control the inflection by a proper
selection of the midpoint knot parameter $t$.

However, the curve shapes are determined not only by the magnitude
of the tangent, but also by the direction of the tangent. The information
about the magnitude of the tangents alone is not sufficient to shape the
curves.

Both the slope vectors $\mathbf{D}_L$ and $\mathbf{D}_R$ can be decomposed into two
components: the component in the direction of the vector $(\mathbf{P}_2 \cdot \mathbf{P}_1)$ and the
component in the direction of the vector $(\mathbf{P}_3 \cdot \mathbf{P}_2)$. Figure 3.8 shows the
decomposition of the right slope vector into individual components. The
vector $(\mathbf{P}_2 \cdot \mathbf{P}_1)$ is translated to point $\mathbf{P}_3$ in order to facilitate the addition of the
two component vectors. The vector $\mathbf{D}_{6R}$ is the component of the right tangent
in the direction of the vector $(\mathbf{P}_2 \cdot \mathbf{P}_1)$ and similarly for the vector $\mathbf{D}_{aR}$. These
vectors are combined to form the right tangent vector $\mathbf{D}_R$.

Therefore, the relationship between the coefficients of the component
vectors $(\mathbf{P}_2 \cdot \mathbf{P}_1)$ and $(\mathbf{P}_3 \cdot \mathbf{P}_2)$ as determined from the endpoint knot parameter
value $t$ gives a radical role in controlling the tangential directions of $C^1$ free-
form splines.

Let the endpoint knot parameters be a function of the magnitudes of
the vectors $(\mathbf{P}_2 \cdot \mathbf{P}_1)$ and $(\mathbf{P}_3 \cdot \mathbf{P}_2)$. That is, the values of the endpoint knot
parameters are constrained such that they are dependent on the distances

between the neighboring points. Then Type 1 spline is defined such that the dependency can be expressed as,

$$\mathbf{P}_i \cdot \mathbf{P}_j = \frac{1}{1+k} \cdot \mathbf{s}$$

$$\mathbf{P}_i \cdot \mathbf{P}_k = \frac{k \cdot b}{1+k} \cdot \mathbf{s} \qquad (3.7)$$

where

vector $\mathbf{s}$ is denote the seat vectors in the direction of the vectors $(\mathbf{P}_i \cdot \mathbf{P}_j)$ and $(\mathbf{P}_i \cdot \mathbf{P}_k)$ (referred to as component vectors, hereafter) respectively.

Dividing the equations, we get

$$\left|\frac{\mathbf{P}_i \cdot \mathbf{P}_j}{\mathbf{P}_i \cdot \mathbf{P}_k}\right| = \left(\frac{1}{k}\right)^2$$

If we define the knot division ratio as $k$ (0.6), then the square of the knot division ratio is equal to the direct length ratio between data points. Therefore, the midpoint knot parameter $k$, and the scaling constant $s$, in terms of the chord length are

$$k = \frac{1}{1 + \sqrt{\left|\frac{\mathbf{P}_i \cdot \mathbf{P}_j}{\mathbf{P}_i \cdot \mathbf{P}_k}\right|}}$$

$$s = \sqrt{|(\mathbf{P}_i \cdot \mathbf{P}_j)(\mathbf{P}_i \cdot \mathbf{P}_k)|} \qquad (3.8)$$

If Equation 3.7 is expanded using Equation 3.8,

$$\mathbf{D}_i = (1-k) (\mathbf{P}_i \cdot \mathbf{P}_j) \frac{1}{|\mathbf{P}_i \cdot \mathbf{P}_j|} (\mathbf{P}_i \cdot \mathbf{P}_j)$$

$$= b(\,s) (s)\, \mathbf{s}$$

$$\mathbf{D}_k = \frac{(1+k)^2}{k} (\mathbf{P}_i \cdot \mathbf{P}_k) \frac{(s)}{s} (\mathbf{P}_i \cdot \mathbf{P}_k)$$

$$- (1 + \lambda + \lambda^2) = 0 \tag{2.9}$$

Hence, the component vectors of the left magnet $D_1$ are of identical magnitude. The same holds true for the right magnet $D_2$, except that they are now sealed by a factor of $(\lambda)^2$.

A geometrical consequence of the above expression for Type I spline is that the tangential direction in the direction of the two unit vectors $a$ and $b$, as shown in Figure 3.4. The property that the tangential directions are identical, is an expected result, since the Type I spline is $G^2$ continuous.

Moreover, the magnitude ratio of the left magnet to the right magnet is proportional to the knot duration ratio $\tau : (1 - \tau)$.

The advantage of the Type I spline can be explained by comparing it with the curve shapes of the ordinary cardinal splines. Figure 3.6 shows the shape vectors of the ordinary cardinal splines. The vectors $D_1$ and $D_2$ denote the left and right tangents at point $P_2$, respectively. Note that the magnitude and direction of the tangents determine the shape of the curve. In the interval $(P_1 P_2)$ and $(P_2 P_3)$, to be noticed, the direction of the vectors $D_1$ and $D_2$ is identically set to the vector $(P_3 P_1)$ in ordinary cardinal splines. In other words, the tangent vector is one half the sum of the two component vectors $(P_2 P_1)$ and $(P_3 P_2)$. Since the vector $(P_3 P_2)$ is greater in magnitude, the corr reach is somewhat the vector $(P_3 P_2)$. Being influenced by the small component vector, the curve shows a band before it reaches point $P_2$. Such bending itself does not cause any problem as long as a smooth transition from point $P_2$.

What matters, and causes the wiggles, is the magnitude of the tangent vectors. Figure 3.10 shows two extreme cases of the ordinary cardinal spline. The curve marked as a) closes the behavior of the ordinary cardinal spline

when point $P_0$ moves away from $P_2$ to the direction of the vector $(P_1-P_2)$. As the distance between the two points $P_0$ and $P_1$ increases, the magnitude of both tangent increases since

$$D_s - D_{s'} = \frac{1}{3}(P_1-P_2)$$

in the ordinary cardinal spline. Vector $D_{s2}$ in figure 5-20 represents three tangents. With point $P_0$ moving farther away, the curve tends to prepare for directional change and bulges in a similar pattern as before, so these $D_s$ show the magnitude of the left tangent while the right tangent of $D_{s2}$ shows, so that the reduction in magnitude of the left tangent makes the curve closer. An inverse case of the reduction is that the left tangent gets longer while the right tangent is kept small, and the spline becomes a pair of line at the point $P_1$.

Figure 5-41 illustrates the curve shape of Type 1 spline for the inner data points as in figure 5-9. The direction of tangent at the endpoint of last unit vectors is to aim in such model that each unit vectors becomes identical with the axis of the reflection. In the left and right unit vectors of Type 1 spline at this endpoint magnitude. The same is true for Type II and Type III splines in its developed shortly. Compare Types 5-11 with figure 5-41 as the curve in figure 5-11 represents the three magnitude of the tangents to flip at the midpoint of the last axis this is reduced while that of the right tangent at the midpoint of the last is approximately the same as before. These magnitudes keep the left and right tangents at the point $P_1$ aplias to reduce direction can be explained mathematically.

The tangents of Type 1 splines in terms of the distance between data points can be derived from Equations 3.8 and 3.9. It follows that

$$|D_1| = \sqrt{3} \cdot x$$

$$|D_2| = \sqrt{2} \cdot (0.6) \cdot x$$

with

$$b.1 = \frac{\sqrt{|P_2 P_2'| \cdot |P_1 P_1'|}}{1 + \sqrt{\frac{|P_2 P_2'|}{|P_1 P_1'|}}}$$

$$(3.42) = \frac{|P_1 P_1'|}{1 + \sqrt{\frac{|P_1 P_1'|}{|P_2 P_2'|}}}$$     (3.10)

As the vector $|P_2 P_2'|$ decreases, the value of b.l will decrease so that the magnitude of the left tangent is reduced. Meanwhile, the magnitude of the right tangent increases as can be verified by examining Equation 3.10. Therefore, the magnitude of the left tangent is inversely proportional to the distance between $P_2$ and $P_2'$, at the Type 1 spline. It is this adaptability of the tangential magnitude that removes the wiggles, which typically appears when a curve in a sheer ascent has to adapt its tangents with constantive departures.

Comparing Figures 3-8 and 3-11 in the interval $[P_3, P_4]$, we notice that the Type 1 spline has slightly more swing in that region. This swing is a result that the magnitude of the right tangent of the Type 1 spline is greater that than of the ordinary cardinal spline due to the tendency of the Type 1 spline to be inflective in the longer interval.

In summary, the distinct nodes of the midpoint knot parameters of Type 1 splines is made to be proportional to the square root of the chord

lengths among data points. Compared with the arbitrary optimal splines, Type 1 splines added improvements in removing the wiggles as the cluster intervals, while showing more swings in the larger intervals. Moreover, the tangential direction is set to be always following the two unit vectors formed by the three consecutive data points.

## Type 2.2 Two Point Splines

Type 2 splines encompass those splines whose endpoint knot parameters have the following relationship with the chord length between data points:

$$P_5 P_6 = Qd_1 + x$$
$$P_9 P_8 = Qd_2 + b$$

(2.12)

where

the vectors is denote the unit vectors in the direction of the vectors $(P_5 P_6)$ and $(P_9 P_8)$ respectively.

In other words, the midpoint knot parameter $t$ is constrained by

$$\left[\frac{P_5 P_6}{P_9 P_8}\right] = \left[\frac{Qd_1}{Qd_2}\right]$$

so that the midpoint knot division ratio is the same as the rate of the chord length. This distinguishes Type 2 splines from the Type 1 splines. Solving Equation 2.11, we get the endpoint knot value $b$ and the scaling constant $t$ as,

$$b = \frac{-P_5 P_8}{|i_1|}$$

$$t = |P_5 P_6| + |P_9 P_8|$$

(2.13)

Substitution of Equation 2.11 into Equation 2.6 yields the expression for the left and the right tangents as,

$$D_x = 0.40 \, (F_1 P_1) \frac{\Sigma^2}{S} \, (F_2 P_2)$$

$$= K \, (1.6) \, mV^2/s$$

$$D_y = \frac{(15)^2}{S} \, (F_2 P_2) \, (F_2 P_2)$$

$$= (>F_1)^n/S \cdot 0.5 m V$$

Notice that the right tangent is again the scalar multiple of the left tangent by the factor of $\left(\frac{N}{S}\right)$ guaranteeing the $G^1$ continuity.

Turning our attention to the relationship between the component vectors making up the left and right tangents, we find that the magnitude ratio of the component vectors depends on the respective parameters. An example of a Type II spline interpolating the same data points as in Figure 3-21 is illustrated in Figure 3-12. The rate of the chord length of the line segment $P_4 P_5$ to that of the line segment $P_3 P_4$ is 1: 3.91 for this specific example, so that $s = .22$ in Equation 3-15.

The consequence of the $k$ value is twofold:

1. It determines the magnitude ratio of the left tangent $D_x$ to the right tangent $D_y$ such that

$$\frac{D_x}{D_y} = \frac{1}{N} = .28$$

2. It determines the magnitude ratio of the directional components. As can be seen in Figure 3-11, the rate of the x vector component to that of the b vector component is,

$$\frac{x \text{ component}}{b \text{ component}} = \frac{|F_2 P_2|}{|F_1 P_1|} = \frac{1.6}{.5} = 3.16$$

so that the tangential direction is close to the x vector, x unit vector in the

direction of $\langle P_4 P_3 \rangle$

The shorter the interval $(P_4 P_3)$, the larger that is a component, and the higher weighting is given to the vector of the shorter interval. This holds true both for the left and right integrals, since the right integral is a scalar multiple of the left integral.

In essence, the sum of the dihedral length determines the multipoint least parameter $k$, and the value of $k$ determines the direction and the magnitude of the tangent.

With respect to the absolute magnitudes of the vectors, converting features is found in Type II splines. Inserting Equation 3.12 into Equation 3.13, we obtain

$$|D_L| = \tfrac{1}{2} |v^2 \sqrt{(v^2+v)+h^2}|$$

$$= |P_4 P_3| \sqrt{(v^2+v)+h^2}$$

$$|D_R| = \tfrac{1}{2} (1/v) |v^2 \sqrt{(v^2+v)+h^2}|$$

$$= |P_4 P_3| \sqrt{(v^2+v)+h^2} \qquad (3.14)$$

If $r$ and $(1/k)$ constitutes the two sides of a right-angled triangle, then the magnitude of the left tangent is the length of the line segment $P_4 P_3$ scaled by the hypotenuse of their triangle.

Compared with Figure 3.11, Figure 3.12 exhibits more features in the interval $(P_4 P_3)$, while it shows more inflection in the second $P_4 P_3$, so the vicinity to the left of point $P_3$, the Type I curve is already changing its direction in anticipation of the component that the tangent should track; the unit tangents of the component vectors of point $P_3$. In contrast, the Type II curve did not change its direction yet. Because the incoming vector $(P_4 P_3)$ is

with a distincure component at the tangent at point $P_3$, the major change of the direction happens only after the curve passes through point $P_3$.

Notice the gradual change of tangential direction at point $P_3$ in figures 3-9, 3-11 and 3-12. The level that the tangent at $P_3$ reflects the direction of the incoming vector $(P_2 P_3)$ is highest in figure 3-12. The wiggles or loops can be regarded in the curve shape where tangential direction undergoes a rapid change in a short interval. In view of the fact that such a sudden change of the tangential direction in the shorter interval is available in Type II splines, it can be said that Type II splines are better than Type I splines for removing wiggles.

In addition to these directional consideration, it is important to note that the magnitude of the left tangent versus right tangent plays another key role in eliminating the wiggles. Writers can explicitly.

$$\left[ \frac{R_3}{R_2} \right]_{\text{ordinary cardinal splines}} = 1$$

For instance, the ratio of the tangential magnitude corresponding to Type II spline in figure 3-12 is

$$\left[ \frac{R_3}{R_2} \right]_{\text{Type I} \ , \ \text{Fig } 3-12} = \frac{3}{0.5} = 15$$

while that of Type I spline in figure 3-11 is

$$\left[ \frac{R_3}{R_2} \right]_{\text{Type I} \ , \ \text{Fig } 3-11} = \frac{3.5}{1} = 3.5$$

The smaller the ratio, the larger the reflection in the longer interval and the smaller the reflection brought in the shorter interval. This can be verified by comparing Figure 3-11 with Figure 3-12 as the thinner interval $[P_2 P_3]$ and the longer interval $[P_0 P_1]$.

In summary, Type II splines dominates the midpoint knot parameter's such that the knot division ratio is the same as the rate of the chord length. The tangent at the midpoint is dependent on the magnitudes of the vectors $[P_2 P_3]$ and $[P_0 P_1]$. In Type II splines, the parameter t axis is the tangential direction of the midpoint such that the direction of the shorter interval is weighted with larger values. The same parameter t was on the tangential magnitude such that the shorter interval has less magnitude. As a result of adjustments at both the direction and magnitude, the material of waggles in Type II splines is more pronounced than in Type I splines.

### Type III Of Free-form Splines

Because we already have defined and analyzed two types of viscosity continuous splines, we will describe the third type of splines so to be embedded in a more general expression. This general expression leads to the definition of a shoot-centric construct-I.

Type III splines are represented by the relation

$$P_2 P_3 = (1-t)^2 + b$$

$$P_4 P_5 = 10t^2 + b$$
  (3.16)

where

b is the endpoint knot parameter;

t is the ability constant;

a, b denote the unit vectors in the direction of the vectors $(P_2 P_3)$ and.

$TP_q$-$E_r$) respectively

the endpoint knot parameter it is constrained by

$$\begin{vmatrix} P_q E_r \\ P_q E_{(q+2)} \end{vmatrix} _{\text{type II}} = \left( \frac{\Delta_q}{\Delta_{q+1}} \right)$$  (3.17)

Compare this equation with the constraints a Type I and Type II splines

$$\begin{vmatrix} P_q E_r \\ P_q E_{(q+2)} \end{vmatrix} _{\text{type I}} = \left( \frac{\Delta_q}{\Delta_{q+1}} \right)^{\frac{1}{2}}$$

$$\begin{vmatrix} P_q E_r \\ P_q E_{(q+2)} \end{vmatrix} _{\text{type II}} = \left( \frac{\Delta_q}{\Delta_{q+1}} \right)$$  (3.18)

The difference in the chord length ratio effects endpoint knot parameter $k$ such that

$$k = \frac{1}{s + \left[ \begin{vmatrix} P_q E_r \\ P_q E_{(q+2)} \end{vmatrix} \right] _s}$$

where:

$s$ = 2 for Type II splines, a chronicalous set of splines can be practiced

$s$ = 1 for Type I splines.

$s$ = 1 for Type II splines.  (3.19)

Note that by varying the $s$ value, a continuous set of splines can be practiced. The curve holds true for the endpoint knot parameter $k$. For instance if it is equal to 0.5, then the characteristic of the spline will be somewhere in the middle of Type I and Type II splines. Aside from the three types of GP splines, an almost unlimited number of splines can be produced. Therefore, we define the constant as the characteristic constant of the GP splines.

As a matter of fact, we have a single degree of freedom in the generation of the curve described by Equation 3.19. As long as we wish to incorporate the chord length information in the shape of a curve the characteristic constant $d$ will alternate the $k$ value, and vice versa.

Depending upon the design purpose, we can obtain a variety of curves which are visually continuous. Provided with a proper control of the characteristic constant, the curve can be made to relieve or to flatten while preserving visual continuity. An example using the characteristic constant, the control of motion path in an interactive design environment, will be developed later.

Subsequent evaluation of the left and right tangents yields the component ratio of

$$\frac{k \text{ component}}{h \text{ component}} = \frac{\partial P(1)/\partial P^T}{}$$

where

Qn,k0    :    $T_l$ — 0 for Type III spline
                      $R_r$ — 1 for Type I spline
                      $U_u$ — 0 for Type II spline                                    (3.20)

The magnitude ratio of the left tangent to the right tangent becomes

$$\left[\frac{R_{ri}}{R_{r,\text{Type III}}}\right] = \left[\frac{R_{ri}}{R_{r,\text{Type I}}}\right] = \left[\frac{R_{ri}}{R_{r,\text{Type II}}}\right] = \frac{1}{4}$$                (3.21)

the left tangent of the curve $\bar{C1}$

the right tangent of the curve $\bar{C1}$

$\bar{C1}$: $z_3 = a + 5$

$\bar{C1}$: $\alpha = 30, t = 67$
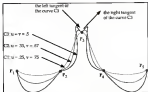
$\bar{C1}$: $z_2 = 30, t = 70$

$F_1$

$F_2$

Figure 5-7 The continuity effect of the $C^0$ free form splines. By controlling the $\alpha$, $t$ parameters such that their sum becomes one, a smooth path enabling each visual effect or an incoming fable can be achieved.

Figure 3-9. Decomposition of the slope vectors $a$, $b$ denote the unit vectors to the direction of the vectors $(P_2, P_c)$ and $(P_c, P_s)$, respectively. $D_{B1}$ denotes the tangential component of the right tangent, to the direction of the $s$ vector. Similarly for $D_{B2}$. In $G^1$ class splines, the magnitudes of the left tangent $D_L$ and the right tangent $D_R$ can be different, but the direction must remain the same.

Figure 5-9 The tangents of an ordinary cardinal spline. Left end tangent and right tangent have the state magnitude and direction.

(a) a curve (knot) appearing
when $\mathbf{r}_0$ moves away from $\mathbf{r}_1$

(b) a curve appearing
when $\mathbf{r}_0$ is drawn near $\mathbf{r}_1$

Figure 3-50: Tangent constraints that cause the wiggles in an ordinary cardinal spline. $\Omega_{C1}$ and $\Omega_{C2}$ represent the tangents associated with the spline C1 and C2, respectively.

Figure 3-11. The integrals of a Type LC[^?] free form spline. a, b denote the unit vectors in the direction of the vectors $(P_2 P_1)$ and $(P_3 P_4)$, respectively.

Figure 3-12. The tangents of a Type II $C^1$ free form spline: $t_i$ denote the unit vectors in the direction of the vectors $(P_1P_2)$ and $(P_2P_3)$, respectively. The ratio of the magnitudes of the a components to the b components is $2.5 : 1$ for this example.

Figure 3-15. Comparison of a Type III G¹ iron farm spline with an ordinary cordinal spline. **x** is the tangent at the side vertex in the direction of the vectors (P₁-P₂) and (P₄-P₃) respectively. The magnitude value of a component is b component is 1 : 3.5 for this example.

### Experiments on the Removal of Margins

In the previous section, three types of $G^2$ free form splines were developed by imposing chord length dependencies on the determination of the multiple knot parameters. In this section, we compare the visual properties of the three types of free form splines in terms of reshaping the weight problem, which was the main motivation of this chapter.

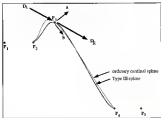A wiggle occurs when the distance between two neighbouring alpha carbon atoms is small compared with its subsequent distance. The variation of the distances appears in the form of tangent constraints at data points. In this respect, all three types of free form splines are eligible for the interpolant of the backbone atoms, since they were developed in such a way that the tangent constraints of ordinary cardinal splines can be alleviated.

However, each type of free form spline has its own characteristic visual performance in figure 3-16, the ordinary cardinal spline method was used to interpolate the data points, $P_j$ through $P_k$. As can be readily seen, the curve shows a wiggle in the interval $[P_i, P_k]$, which is a relatively short interval compared with the interval $[P_j, P_i]$.

Figures 3-15 through 3-17 show the interpolation of the same data points using Type II, Type I, Type III free form spline, respectively. All of them contribute to eliminate the wiggle in the interval $[P_i, P_k]$, thus making them adequate for the computation curve for the backbone curve

Figure 3-18 shows a closer view of the behaviour of the three types of free form splines in the region $[P_i, P_k]$. The Type II spline has the lowest level of constraint in the region, while Type I and Type III show an increased level of constraint. However, the curvature at this for the tangent interval $[P_i, P_k]$

Figure 3-31 illustrates the tonal properties of the curves starting from point $P_5$ up to about half way to point $P_3$. Type I and Type III pillars situate predict the curve of the ordinary nonlinear spline, while the Type II spline exhibits a considerable amount of overshoot.

Therefore, we can conclude that the level of overshoot in the shear interval is proportional to the radius of characteristic constant $k$ in Equation 3.19, while that of the longer interval is inversely proportional to the characteristic constant.

The decision as to which type of line loses spline is better depends on the various design purposes. If a subjective matter. A designer may regard the swing in the longer interval as a mere natural one, or he might prefer a straight line in the longer interval and trade off the swing with that of the shorter interval.

Each type of line loses spline can be extended to a certain interpolation scheme. For instance, the bilinear blending method can easily connect a combination of space curves into three-dimensional surface shaper. Figures 3-31 and Figure 3-32 show Curial patch (Coon9, Part1) applied to different interpolating schemes. The surfaces show the same trend as was seen in the case of a space curve. The surface blended by a nonlinal spline tends to be bent at the small patch on the left side, while in the half patch produced by a Type II spline shares some influence on the half patch on the right side. Side views (see Figures 3-33 and 3-34) demonstrate this property more clearly.

Figure 3-11. A wiggle produced by an arbitrary cardinal spline.



Figure 3-13. Removal of the wiggle by a Type II G² low form spline.

Figure 5-16  Removal of the wiggle by a Type I G7 free form spline.



Figure 5-17  Removal of the wiggle by a Type II G7 free form spline.

Figure 3-18. Comparison of two form splines in the shorter interval $[P_2, P_3]$



Figure 3-19. Comparison of two form splines in the longer interval $[P_3, P_4]$

Figure 3-37. Sample point locations used for the construction of surface patches in Figures 3-41 through 3-51.
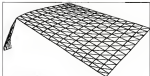


Figure 3-51. Coons' patch produced by an ordinary cardinal spline.

Figure 3-22. Corner patch produced by a Type II spline.



Figure 3-23. A side view of Figure 3-22.



Figure 3-24. A side view of Figure 3-22.

## Comparison with Truncated Cardinal Splines

Having accomplished the task of removal of singular in the last section, we devote this section to the comparison of the properties of the $G^*$ free Euler splines with the well-known truncated cardinal splines [brand]. In this section, the Hermitian form of our free Euler splines will be compared with that of the truncated cardinal splines. To distinguish free Euler splines from the truncated cardinal splines, the superscript for truncated cardinal splines will be derived in the context of direct length parameterization. This way, the difference between the truncated cardinal splines and our free Euler splines making use of the closed length orthonormatrices can be clarified. The visual properties of the two approaches will be explained in terms of removal and addition of tension.

A cubic spline can be expressed in terms of directional vectors instead of positional vectors (for location of data points). By a simple manipulation, the point tensions can be converted into directional vectors. For instance, the general motion expression of a cubic spline is

$$P(s) = U \cdot A \cdot [P_1 \; P_2 \; P_3 \; P_4]^T$$

Instead of writing the point matrix as $[P_1 \; P_2 \; P_3 \; P_4]^T$, one can rewrite it in terms of directional vectors. In the method, instead of $P_2$ and $P_3$ are introduced or point $P_2$ and $P_4$, so that the point matrix becomes $[P_1 \; P_2 \; (P_2) \; P_4]^T$.

Consequently, the point matrix is made up of endpoints $P_1$ and $P_4$ on both ends of the curve, and two vectors $(P_2 P_2)$ and $(P_3 P_3)$ representing the slopes at $P_2$ and $P_3$, respectively. As a matter of fact, piecewise Hermitian interpolation allows for slope vectors other than $(P_2 P_2)$ and $(P_3 P_3)$, and the

representation belongs to a special case of the general Hermite manipulation regression. Using the directional notation, the two form syntax in Hermite form can be written as

$$p(x) = 1/z \; A + z \; 1/z \; A \; [P_1 \;\; P_2 \;\; P_3 \;\; P_4] \; D_1 P_2 D_1^T$$

with

$$
A = \begin{bmatrix}
\frac{1}{1+z} & \frac{2}{1+z} & \frac{z}{1+z} & D + \frac{1}{1+z} & r(\text{1+z}) & \frac{1}{1+z} \\[6pt]
\frac{2}{1+z} & \frac{1}{1+z} & \frac{1}{1+z} & D + r & \frac{1}{1+z} & (1+r) \frac{1}{1+z} \\[6pt]
\frac{1}{1+z} & \frac{1}{1+z} & z & & & \\[6pt]
\frac{1}{z} & 0 & & & & \\[6pt]
1 & & & & & 
\end{bmatrix} \tag{3.32}
$$

Notice that most of the elements in each column of matrix $A$ are fractions of the $x$ and $r$ parameters. Thus the values of the matrix elements will reflect the change of the endpoint knot parameters $x$ and $r$, if the parameters are associated with the chord length between the data points. However, as will be explained next, only the last two columns of the matrix $A$ will be dependent on chord length if the tensioned cardinal splines are maintained in Hermite form.

The use of tensioned cardinal spline has been a conventional method of adjusting curve shapes, using a parameter called tension, while preserving the $C^1$ continuity. Ordinary cardinal splines can be expressed in terms of the Weier central points (WcP) which inside axe of the knot parameters to attain the tangent continuation at joints. As a result, the cardinal spline can be written as a function of knot parameters. An additional constraint then the

kost applying be proportioned to the chord length leads to the tensioned cardioid splines.

In figure 5.20, data points $P_i$ through $P_n$ are to be approximated by a tensioned cardioid spline, and they are assumed to be data points corresponding to the parameters $u$ which takes the values of $u_i$ through $u_n$.

The tensioned cardioid spline can be derived starting with the construction of Bézier control points. If points $P_j$ and $P_k$ correspond to Bézier control points $b_0$ and $b_3$ respectively, then points $b_1$ and $b_2$ are set called the post Bézier control points positioned in relation to the post Bézier guide points $b_1$ and $b_2$, respectively. The position of these control points can be evaluated by defining the knot space $u_k$ as

$$\Delta_i = u_{i+1} - u_i,$$

Using a method for tangent estimation known as FMILL [Far65], two inner Bézier points can be expressed as

$$b_1 = \frac{\Delta_i}{3(\Delta_{i-1}+\Delta_i)} P_j$$

$$b_2 = \frac{\Delta_i}{3(\Delta_i+\Delta_{i+1})} P_k$$

where

$$\hat{P}_j = \frac{P_j}{\|\Delta_{i-1}\|}$$

$$\hat{P}_k = \frac{P_k}{\|\Delta_i\|}$$

with

$\|u\|$ standing for the norm of vector $u$.

Using these values, the general Hermite expression of a cardioid spline in the interval $[P_j, P_k]$ becomes,

$$P(u) = P_j f_1(u) + c_j f_2(u) ...$$

$$+ \frac{\Delta_2}{3(\dot{x}_1+\dot{x}_2)} \left[ x (2x^2-6\bar{z}^2+3z_1z_2+3z_1z_3+3z_2z_3) \right.$$

$$\left. - \frac{3}{2}\bar{z} \,x^2 + z \right] 2x \left[ N \,(F_x \,F_y \,)F_x F_y \,\bar{x}_y \bar{y}_y x^2 \right]$$

with

$$N = \begin{bmatrix}
2 & -2 - \frac{2\bar{z}_2}{(\bar{x}_2 \bar{F}_y)(\dot{x}_1+\dot{x}_2)} & \frac{2\bar{z}_3}{(\bar{x}_1 \bar{F}_y)(\dot{x}_1+\dot{x}_2)} \\
-3 & 3+\frac{6\bar{z}_2}{(\bar{x}_1 \bar{F}_y)(\dot{x}_1+\dot{x}_2)} & \frac{-6\bar{z}_3}{(\bar{x}_1 \bar{F}_y)(\dot{x}_1+\dot{x}_2)} \\
0 & 0 & 0 \\
1 & -1 - \frac{4\bar{z}_2}{(\bar{x}_1 \bar{F}_y)(\dot{x}_1+\dot{x}_2)} & \frac{4\bar{z}_3}{(\bar{x}_1 \bar{F}_y)(\dot{x}_1+\dot{x}_2)}
\end{bmatrix}$$

$$(III.91)$$

Notice that the first two columns of matrix N are the same as columns of matrix N with the Hermite form of the free loco spline. The anchory δ₂ in the sense that it is last transitional crotchet spline represents one instance of this expression where

$$\dot{x}_2 = \dot{x}_3 = \dot{x}_4 = 1$$

$$\dot{x}_1 (F_y F_y \dot{x}_1 \dot{x}_1 F_y \dot{x}_1) = 1.$$

The loco space variables δ₂ and δ₃ can be parametrized by chord length such that

$$\frac{\Delta_2}{\Delta_1} = \frac{|(F_y \bar{F}_1)|}{|(F_y \bar{F}_2)|}$$

$$\frac{\Delta_3}{\Delta_1} = \frac{|(F_y \bar{F}_3)|}{|(F_y \bar{F}_2)|}$$

This interpolation, called 'chord length parameterization', yields a $C^2$ continuous crotchet spline so that the derivative with respect to the parameter $u$ is continuous

We can prove that chord length parametrization is $C^1$ continuous as the following way. Let

$$s = \frac{k_1}{(\mathbf{P}_1 \mathbf{P}_2)^{1/2}(k_2+k_3)}, \text{ and}$$

$$\mathbf{h} = \frac{k_1}{(\mathbf{P}_1 \mathbf{P}_3)^{1/2}(k_2+k_3)}$$

Then the above Hermite form of cardinal spline reduces to

$$P(u) = U \cdot \mathbf{H} \cdot [\mathbf{s} \ \mathbf{P}_1 \ \mathbf{P}_2 \ \mathbf{P}_3]^T$$

with

$$\mathbf{H} = \begin{bmatrix} -s & 2-h & s-2 & h \\ 2s & h-3 & 3-2s & -h \\ -s & 0 & s & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} \tag{5.26}$$

which is exactly the recommended tensioned cardinal spline. Notice that the routine agrees with that at the four form spline, if both midpoint knot parameters are equal to $\frac{1}{2}$. Therefore, the left tangent at point $\mathbf{P}_2$ is

$$\frac{dP(u)}{du}\bigg|_{u=1} = h(\mathbf{P}_3 - \mathbf{P}_2),$$

and the right tangent evaluated by a curve Q(u) in the right of point $\mathbf{P}_2$ is,

$$\frac{dQ(u)}{du}\bigg|_{u=0} = h(\mathbf{P}_3 - \mathbf{P}_2). \tag{5.25}$$

In order to assure $C^1$ continuity at the point $\mathbf{P}_2$, the value of $s$ and $h$ must be the same. This is the reason why the tension parameter $t$ replaces both $s$ and $h$.

as the terminated cardinal spline, thus removing the main degree of freedom of varying the two independent knot spacings.

Although the locality of a spline is determined by the number of changed curve segments with a modification of a data point position, the locality can also be described by the number of changed curve segments with a change at a parameter value. If the shown parameter '+' to 'b + 4' or to have a different value, then the same value of '+' should be applied to entire curve segments. Otherwise, the joint between two curve segments whose tension parameter differs, fails to be $C^2$ continuous. Consequently, the same tension value is applied reductionantibody to all curve segments, regardless of the distance between the segments.

Figure 3-30 illustrates the use of a terminated cardinal spline for the removal of a wiggle in the curve caused by an ill-fitting tension. In the interval $[P_a P_b]$ the wiggle is removed. Nevertheless, the interpolation in the region $[P_2 P_3]$ is almost linear, since the same tension value used to alleviate the wiggle was applied in the region with constant because of the $C^2$ continuity constraint. Moreover, the long unravel length in $[P_2 P_3]$ acts to flatten the curve. Besides being unwanted, each linearity problem was a small amount of deflection on the curve, which is one of the reasons why the wiggle appear was employed.

Figure 3-31 illustrates another appeal of terminated cardinal splines with correct corresponding to two different tension values. A constant tension value of .5 coincides with an ordinary cardinal spline, which is not terminated. In a curve with a tension of 1.6, the left and right tangents at point $P_a$ are denoted by $D_b P_a$ respectively. Suppose we are not satisfied with the tension of the curve with a tension of .5 in the interval $[P_2 P_3]$. Then we might be loosening the degree of inflection in this region by raising the tension to 1.0.

However, the problematic area of the curve at hand is the region $(P_b, P_f)$ that runs because of the application of the same section value in the short interval, the curve shows overshoot in this region. In addition, the overshoot of the temporal magnitude thus the left and right tangents must fix the same makes the curve show greater curvature (and point $P_f$, thus near point $P_b$. Therefore, installing these shapes is hardly acceptable as a reasonable interpolation.

In addition, the endpoint least parameters of $G^2$ (two knots spline can vary for each data point. A parameter value at a data point affects only two curve segments adjacent to the point so that the parameter can be applied locally. More importantly, distance collocation can be incorporated into the spline in such a way that could adjust tension in the shaper as shown in interval. For instance, while enlarging the weight at the shorter interval, the Type II spline can introduce an arbitrary amount of inflection at the linear region as were shown previously in Figure 3-13. Therefore, $G^2$ that form splines can be said to have more flexibility compared with the restrained cardinal spline.

It is important to note that the random assignment of a and b values can produce $G^2$ continuity at some data points, or can be deduced from Equations 3.35, since the left tangent is a scalar multiple of the right tangent. Nevertheless, these random values affect the curve shapes between the intervals in such an unpredictable way as to be unacceptable as a proper interpolation scheme. The curve may exhibit extreme oscillation or inflection between the endpoints, and the major factor that causes this behavior is that the curve should meet the very restrictive constraint that the tangential directions must be $(P_b, P_f)$. Such unacceptable curve behavior between the

midpoints has left the tensioned confined sphere with even restrictive $C^0$ continuity.

Although the Eulerian lines of confined sphere appears to embody the total spicing parameter, it is different from just formulations at free form spheres. The first two columns of the R metros of the vertical Hermite form fit Equation 3.25 are constant, even in chord length parameterization. The only way the knot spicing could alter the spline is by way of the last two columns of the R matrix. Although a modification of these two columns subsequently alters the magnitude of the slope at points $P_0$ and $P_1$, namely ($P_0 P_1$) and ($P_1 P_0$), only the magnitude of these vectors can be changed as a result of the variation in the knot spicing.

In contrast, free form splines do not exhibit unknown of Equation 3.27 to influence the curve shape, so that the resulting curve is not strictly a change of integral magnitude, but the adjustment of the tangential directions. To make a curve to $G^1$ continuous, there should be a typical of freedom in selecting the direction of tangents of given data points, since the directions nontrivially affects the nature of a curve IdentiG. Otherwise, the curve will undergo an unwanted oscillation between these data points to adjust itself to the fixed tangential directions at these points.

One approach to remove this unevenness or the tangential direction of the data points is presented by Kochanek (Kochk) bistarg with the Hermite form of arbitrary cardinal splines: he has modified the slope component of the point matrix, in such a way that the tangent direction can be adjusted to our terms. Kochanek's decomposed the fixed tangential direction at point $P_i$ into two components,

$$P_0 P_0 = D_0 P_0 + (P_1 P_0)$$

In order to control the curve shape, these two vectors are weighted by certain parameters such that

$$\mathbf{P}_4 \, \mathbf{P}_3 \;=\; \delta_3 \, (\mathbf{P}_4 \mathbf{P}_3 \mathbf{P}_1 \wedge \mathbf{P}_5 \mathbf{P}_3)$$

where the weighting factors $\delta_2$, $\delta_3$ are freely modified to control the tangent and, thus the shape of the curve. In a general expression, his method in Hermite form corresponds to

$$\mathbf{P}(u) \;=\; (u^3 \; u^2 \; u \; 1)
\begin{pmatrix}
2 & -2 & 1 & 1 \\
-3 & 3 & -2 & -1 \\
0 & 0 & 1 & 0 \\
1 & 0 & 0 & 0
\end{pmatrix}
\begin{pmatrix}
\mathbf{P}_i \\
\mathbf{P}_j \\
\mathbf{D}_i \\
\mathbf{D}_j
\end{pmatrix}$$

with

$$\mathbf{D}_i \;=\; \tfrac{1}{2}\,\delta_3 \,(\mathbf{P}_{i+1}\mathbf{P}_i \wedge \mathbf{P}_{i-1}\mathbf{P}_i)$$

Although he extended the same scheme to the point that the left and right tangents are not the same, our discussion in that section will be restricted to the identical tangent case. Further dimension will follow at the next section.

This method can be characterized by two observations. First, the curve is not generally $G^2$ continuous, since the left tangent is identical to the right tangent. Rather, it can be said to be $G^1$ continuous where the left and right tangents are the same. Therefore, the variety of the curve shapes extractable from the $G^1$ class could not be covered. Second, despite the fact that a $G^1$ continuous, the curvature is not based on the chord length information. Therefore, the method works well when the data points are equidistant, but it demands iterative readjustment of the weighting factors if the data points are not equally spaced. In practice, situations when data points are equidistant are

ring. The difficulties involved in the iterative manipulation of the slope
control parameters are mentioned in the next section.

In summary, two basic splines are more flexible than unmatched
cardinal splines in that the direction of tangents can be freely adjusted. Even a
mathematical viewpoint, they belong to the class of $C^1$ spline encompassing
more diverse curve shapes than $C^2$ splines, which are a special case of $C^1$
splines. Moreover, the chord-length dependency of the basic splines
facilitates the manipulation of general random-spaced date points which are
not necessarily equidistant. The flexibility in selecting the tangential direction
of free form splines will be exploited further in the next section, in
conjunction with a method to control the curve shape interactively.

## Application: An Intermittent Motion Path Generator

This section extends the use of free form splines to an intermittent
animation tool to generate a motion path. More frequently, animators design
a motion path iteratively, rather than resorting to the three types of
automated motion described in previous sections.

To date, research into control of curve shapes has mainly relied on the
introduction of new parameters [Nutt5, Barr85, Bant84, Pickl, Kochl]. The
designers had a tough idea of what the curve will look like when they change
the parameters. Nevertheless, the exact behavior of the curve would hardly be
predicted until it was left to the designer to tentatively vary and fine tune the
parameter values until they reach a satisfactory curve shape. This happens
because the parameter values are entities which cannot directly interact with
our visual perception, while the data points related to the parameters are
geometric. For instance, one may question exactly how a basis spline will

believe that a does point if the blue parameter of the point is charged from $\beta$ to $\gamma$.

The greatest advantage of $G^2$ flux form utilizes developed center lies in the fact that the specification of surgery vectors at data points can directly lead to the piecewise formulation of the splines. Exploiting this advantage, the only thing designers need to do is to draw tangential directions for individual data points. They can specify how the curve will proceed at the data points simply by drawing a vector at that point; it is the responsibility of $G^2$ flux form splines to produce curves which exactly match these tangents. This way the correlation between computer and designer becomes easier and simpler because it is no longer necessary to predict the visual effect that the under parameter will cause.

This section reviews previous research from the control of curve shapes and presents a new method enabling more flexible control of the curve shapes based on $G^2$ flux form splines.

### Relevant Works

Before investigating previous works, let us briefly examine the dual usage of a spline curve.

A spline can be used for two purposes in interaction. Figure 7.10 illustrates this point. At time $t = t_1$, a snapshot at data point $P_1$, $P_2$, $P_3$, $P_4$ can be interpolated by using a spline to produce a spline curve. Suppose the four points change their positions as a function of time, such that their locations appear as in Figure 7.2b at time $t = 2$ and $t = 3$. With respect to point $P_1$, three data points are given to the data sets. By assuming the three locations for the same point $P_1$, corresponds to generating a locus of

movement to the time set for the data point, and a spline was used for the generation of motion path for individual data points in space. Therefore, control of the shape of the motion path can directly apply to control of the shape in space.

There have been efforts to control the shape of three-dimensional space curves in terms of loci and tension [Koch_Bart01, Bart99_Fo08, Kin99]. The term control of tension represents making a curve tighter or looser at data points [Wit95, Koch6] or on the curve segments between data points [Bart01, Fo08, Si9601]. In general, a curve can be said to be highly tensioned if the curve shape tends to be linear between data points. For instance, a Type II spline exhibits lowest tension relative to other splines, as was shown in Figure 3.19

The term line has been used to represent the variation of tangential direction at this point. The line can be used to visualize the traditional animation effect of following through after an action, or anticipating a movement [Thu01, La99]. Anton hardly came to a sudden and complete stop, but we generally control just their termination point (for example, a hand will follow through after throwing a ball. This tendency to preserve the previous direction can be shown near point $P_2$ in Figure 3-11. The tangent at point $P_2$ resembles the vector ($P_1$ $P_2$) so that the direction of the curve around even after it passes through point $P_2$.

On the other hand, anticipation prepares for a movement. It is a technique to reach the audience's eyes, to prepare them for the next movement and to then expect it before it actually occurs. For example, the tangent at point $P_2$ in Figure 3-10 resembles the vector ($P_2$ $P_3$) rather than the vector ($P_1$ $P_2$). Therefore, the motion path can be controlled to change its

direction before it actually another point $P_o$, and prepares for the next movement

As approximating splines, Reis splines [Rue68, Rue88] were modeled with the inherent capability of controlling the shapes using a pair of scalar parameters, bias and tension. Curve shapes are determined by the specification of two independent parameters. With a fixed bias value, the effect of the variation of tension can be easily recognized. Similarly, a bias effect can be observed only if the tension is not fixed. However, it is hard to predict a curve shape produced by the combination of the two parameters, since a tension between data points can offset the bias on the same data points. In practice, the two parameters interact with each other during the determination of the curves, even if designers think they are independent.

In the area of interpolating splines, Nielson developed the $\nu$-spline [Nie74, Fol90] as a polynomial alternative to Schweikert's exponential-based $\tau$-spline [Sch66], the spline under tension. He insisted that the exponential function is a major requirement in efficient evaluation of the splines and replaced it with a $C^2$-continuous [Bar88] polynomial, where the difference of the second-order left and right derivatives is set as a multiple of the first-order derivative at a data point. The proportionality concern, being set as the tension parameter at the point, liberates the curve as it approaches infinity. Unfortunately, the method lacks a mathematical proof of the existence of a unique solution [Sch83, Sar88]. Experiments on negative tension values in the $\nu$-spline can be found in [Fol90]. However, finding a good tension value is not a trivial matter and no consistent method for doing so is yet known.

Given data points, and the corresponding tangents, Hermite interpolation [Fol90] can be used to produce $C^1$ cubic splines. The tangents should be specified in both magnitude and direction. The tangent direction

can be perceived in such a way that it could fit the desired shapes. However, from a designer's viewpoint, the specification of imperatival amplitude belongs to a difficult task. For instance, one might have no idea how the curve will change if the magnitude of a data point is multiplied by a factor of ten. Moreover, the curve generated by such repetitive manipulation belongs to the $C^1$ class, which has highly restrictive tangential constraints, and has less variety than the $C^2$ class line turn splines.

Rechtsmeck [Recht] has defined both bias and tension as a weighting function of two vectors denoted by subtracting three adjacent data points. He even allowed the left tangent to be different from the right tangent at data points to define a continuity parameter. The subsequent curve shape is defined by the multiplications of the three parameters—bias, tension and continuity. One of the difficulties with this approach is that the combination of parameters may yield unexpected shapes. For instance, specification of zero continuity yields the same curve as a curve with a tension value of one. If an animator decides to set the tension by one, he should strive to make the continuity parameter equal zero to maintain a consistency to his specification of parameters. That is, he needs to have an a priori knowledge of what parameters conflict with one another. Another disadvantage of this approach is that the curves so generated are not only discontinuities in tangential magnitude, but also discontinuities in the direction of tangents. Generally, the curves are neither $C^1$ continuous nor $G^1$ continuous.

To characterize previous research done in the control of curve shapes, it can be said that all of the methods incorporate too many burdens to be useful to designers. For instance, not only one chooses the values of the three parameters—bias, tension and continuity. The animator should be prepared with some knowledge of what bias, tension, and continuity are

Thus he may increase and decrease these values based on that knowledge heuristics he may encounter unexpected curves caused by interference between the parameters, and have to replace the parameter values with what he may think would have a high probability of yielding the desired curve.

## Target Specification Method

From an assistance point of view, the previous method demands too much imagination. These methods require that the operator supply exact values of parameters and then order values could not appear in the usual prediction at the curve shapes. Imagination must be exploited to predict curve shapes from given noble parameters. When these parameters are involved, such visualization is too easy, even though the concept of individual parameters and the possible associations among them are kept in mind.

Central C? too hot igloos allow us to reverse the above procedural direction. That is, an imaginary curve will interpret the inversion values. Imagine a curve interpolating given data points without any concern or knowledge of the concept of the shape control parameters. Thus one can specify the important characters by drawing a line segment starting from individual data points. These longevity is tasks will be used in the formulation of the parameter C? has been option. In this method, it is the programmer's responsibility to abstract the target information into corresponding d parameters of the C? has been option. Therefore, the only task left to the selector is to chose longevity directions based on his imagination. It is important to note that shape control is now done by visual vector reacting instead of the specification of numbers.

$$a = B_1 \cdot P_1 - D_1$$
$$b = B_1 \cdot P_2 - D_1 \tag{3.24}$$

The curve $ba/dx$ near for the tangent $D_a$ at point $P_a$ except that the incoming and outgoing vectors are replaced by $(P_a \cdot P_1)$ and $(P_a \cdot P_2)$, respectively. A direct consequence of Equation 3.9 is that the sums of these two magnitudes is a function of midpoint knot parameter $k$ and the length of the incoming and outgoing vectors.

$$\frac{k}{\sqrt{k}} = \frac{(1-k)^2}{k^2} \cdot \frac{|P_a \cdot P_1|}{|P_a \cdot P_2|} \tag{3.25}$$

Note that the expression is valid both for the left and right segments with different magnitudes. The value of the midpoint knot parameter $k$ is then

$$k = \frac{1}{1 + \sqrt[3]{\dfrac{|P_a \cdot P_2|}{|P_a \cdot P_1|}}} \tag{3.26}$$

Therefore, a unique midpoint knot parameter value is determined given the direction of a tangent. The range of $k$ is confined to $0 < k \leq 1$ by the above equation. In addition, since the $k$ value of one is not defined by Equation 3.5 and 3.6, such a situation can be avoided by assigning an infinitesimally small value. In case a component is zero, which happens when the tangent is perpendicular to the vector $(P_a \cdot P_2)$ in the special case of

$$\frac{k}{\sqrt{k}} = 2.$$

Equation 3.26 reduces to a Type I spline as expected. As can be seen Figure 3-51, the resulting curve is $G^1$ continuous while being faithful to its original tangent specification given in Figure 3-50.

Note that the decomposition of vectors and thus the calculation of endpoint time value it can be handled completely by an algorithm so that endpoints do not have to be constructed since it

Figure 3-33 illustrates another aspect of the tangent specification method from the point of view of tension control. In this figure, the tangents of various segments in two consecutive inflection points, tension control point 1 and inflection point 2 at the region $(P_1 P_3)$. Compare this with the curve in figure 3-31, which shows a single inflection point after point $P_3$. The reflection point 2 is created by making the tangent at point $P_2$ similar to the vector $(P_3 P_4)$ rather than $(P_1 P_2)$. The curve shown inflection before it reaches point $P_4$, to meet the requirement at the tangent direction at that point. Thereafter, the tension can be controlled near both ends of the interval $(P_3 P_4)$, either individually or simultaneously. Often, preference is given to tighter curves to prevent an overshoot. Figure 3-35 illustrates another case of the tangent specification with corresponding curve shapes. In this figure, both tangents at points $P_3$ and $P_4$ are bending toward the vector $(P_3 P_4)$, making the curve beneficial in the interval $(P_3 P_4)$. The more such tangent overrules the vector, the tighter the curve will be constructed.

In $G^1$ two turn splines, a tangent knot parameter divides the curve segment into two intervals centered at that point. Therefore, a value curve in one interval is affected by two adjacent knot parameters corresponding to the two endpoints at both ends of the interval. In addition, two tension points neighboring from one affect the curve segment. Hence modification of a data portion affects four curve segments enclosing the curves not neighboring the point. Compared with this, modification of a tangent at a point affects only two curve segments: the segments to the left and right of the point, whose tensyl properties upon change of the tangent direction can

be predicted tentatively. Therefore, the stagnant specification method facilitates slightly less lateral control of nerve shape without further complicating nearby chron argument.

Figure 3-34 illustrates shape control by a general line form spline which is not $G^2$ continuous. For instance, or instance may want to make the curve in the internal $[P_1, P_4]$ highly tentatived, yet may not be satisfied with the inflection in the internal $[P_4, P_2]$ appearing in figure 3-35. Since the inflection to the left of point $P_4$ in figure 3-35 is caused by the $G^2$ constraint, one line has to sacrifice this continuity to remove the inflection. The $G^2$ continuity has to be traded off for more freedom in controlling the shape. In the figure, $P_{2L}$ and $P_{2R}$ represent the left and right tangents at point $P_2$, while $D_{1L}$ and $D_{1R}$ are those of magnitude $B_{2L}$ and $D_{2R}$ at the same in the internal $[P_4, P_2]$ and $[P_2, P_3]$. When magnitude $B_{2L}$ and $D_{2R}$ and the magnitude $D_{1R}$ make the curve highly tentatived at the internal $[P_4, P_2]$ and $[P_2, P_3]$. The lines caused by magnitude $B_{2L}$ and $D_{2R}$ remove a distortion at the region $[P_4, P_2]$. The manner of the tangents agree with our original inflection at the front spline, when the parameters $s$ and $s$ used remain different.

For instance, the parameter $s$ can be replaced by the value matching the tangent $P_{2L}$ and the parameter $s$ can be replaced by the $k$ value matching the tangent $D_{1R}$ to remove the curve in the internal $[P_4, P_2]$. The recipient least parameter $k$ can be calculated using the same decomposition procedure as described above, and internal case function 3-3 to yield a general expression for line data splines. A comparison of this curve near point $P_2$ and $P_3$, it is shown that the curve near point $P_2$ has more parameters in the $G^2$ line form spline than the one near point $P_3$. If the left and right tangents are more alike at these directions, the generated curve exhibits more uniform the line form data point. Replacing this observation, we can exert control the shape of

discontinuity by adjusting the direction of the left and right tangents of individual data points.

In summary, we have developed an intuitive method for control of the motion path during correction. Compared with the previous approach that is purely number oriented, our tangent specification method lends directly to formulation of parametric curves that match the specified tangent directions. An intentional violation of the visual continuity can also be incorporated into the method where the degree of discontinuity can be controlled by the curvature between the left and right tangents. Moreover, bias, tension, and continuity parameters are integrated into the directional tangents in our method.

## Summary

By stating the characteristics of the visually continuous class of splines, we could remove the undesirable feature of the general interpolating splines, called wiggles.

Our contribution in this chapter is that the visually continuous class of splines has been mathematically formulated in an analytic metric form. To date, no known solution to the visually continuous class of splines has been presented in its analytic form. By extending the class of splines with piece-wise solution, we could define characteristic constant $d$ that could be effectively utilized to control the curve shapes, and to remove the wiggles as well. Finally, as an application of the use of the characteristic constant, an interactive design tool for motion path generation was developed.

knot parameters corresponding to individual points $t_i$

Figure 3-25. Ordinary cardinal spline in terms of the Bézier control points. The knot parameter is shown as mapped into the index curve in the region $P_5 P_6$.



Figure 3-26. A tensioned cardinal spline used to remove a wiggle. The tension value is .85.

Figure 5-27. Slope vectors of truncated cardinal splines $D$ and $D_2$ represent the left and right tangents at point $P_2$, respectively.

Figure 5.20: A spline used for two purposes—for the interpolation of data points in $x$ state-space, and for the interpolation of a force of a data point in the time domain.

Figure 5-29 Direction of tangents specified by an animator

Figure 5.10 Construction of the basis-level parameters by a decomposition of the directional vectors in Figure 5.10 Computed λ values are 5, 25, 5, 25 et ponto P₃ through P₅ sequentially.

Figure 3-30  Final G² free form curve shape matching the tangent specification in Figure 3-29



Figure 3-31  Reduction of tension in both the right of point $P_2$ and the left of point $P_3$. Computed $b$ values are 5, 20, 66, 3 at points $P_2$ through $P_5$ respectively

Figure 5-51 Increase of the tension in the region $[P_1,P_2]$, $D_1$ and $D_2$ represent the tangents at points $P_1$ and $P_2$ respectively. Compared $k$ values are $A$, $B$, $S$, $b$ in points $P_1$ through $D_2$ respectively.

Figure 5-9d: Breakup of $G^1$ continuity. $D_{14}$ and $D_{41}$ are the left and right tangents at point $P_4$, while $D_{1,9}$ and $D_{9,1}$ are those at point $P_9$. Compared it values are -1 for $D_{14}$, 35 for $D_{9,1}$, 8 for $D_{1,9}$ 1 for $D_{41}$. Additional tangent condition is that the right tangent of point $P_4$ has $\theta$ value of 1 and the left tangent of point $P_9$ has that of -1.

# CHAPTER 6
## CONTROL OF MOTION SPEED IN ANIMATION

### Introduction

The display of animation sequences is produced by showing a series of changing still pictures at a constant repetition rate. For video-compatible systems, the frame rate is 30 per second to create the illusion of reasonably smooth motion. In particular, the speed of animation is controlled by the number of frames for a given sequence. For instance, consider the animation of an atom moving through a pair of fixed positions, A and B in space. If the motion between the two postures is captured and displayed with 100 frames this motion will appear about ten times slower than the one with five frames. Replay during system breadcrumb applies the actual effect for direct inspection of movement.

If the motion path of the atom is expressed in terms of a normalized parametric curve P(s), then it can be said that the atom moves from A to B while the parameter s runs from zero to one. That is, the locus of the atom is s dropping from the a parameter domain into three-dimensional space for ease ofview. Frames are required to describe the motion from A to B exclusively, the parameters s may be divided by the, and then the frame is generated by an s value of zero and the second frame is generated by an s value of 1, and the third by an s value of .2, and so on. Thus the question is, "Does the atom exhibit constant motion speed with the atom governed by those parameter values ?" Theoretically, if the distance the ball travels

132

between frames is identical through the entire sequence of frames, the motion speed is correct.

Unfortunately, content splicing in the parameter domain does not guarantee content splicing in the space domain except in the case of a linear motion path. That may be considered as a sampling problem in the parameter space: in general, dense sampling segments the detail of the trajectory, providing a close-texture version of the action, and coarse sampling compresses the detail, providing a time-lapse version of the action. Nevertheless that is not necessarily true for the motion path generated by a cubic spline where a small increment of the parameter u may produce a large distance gap and vice versa. A further requirement concerning the control of the motion speed is "Can we even introduce and deintroduce the motion speed as the parameter moves by proper sampling of the u parameter?" This problem is known as kinematic adjustment or control of motion speed, or motion dynamics in animation, and this chapter is devoted to solving this problem.

In an animation of passive folding, the motion speed of a single atom can be said to determine the speed of the transition from one folding state to another. Error can sometimes explain the passive folding as based on the lilleformming of key frame representing the folding interactions, the control of speed in the interleaved frames is important, figure 6-1 illustrates the case when no motion speed control is taken. Consider the portion of alpha carbon of three sequential data stations $0 = 0$, $0 = 1/2$, and $u = 1$. The alpha carbon at each instant can be assumed to be part of the whole backbone chain representing a folding state. Given the three folding states corresponding to each new instance, spline interpolation can be used to interpolate the motion path of the atom. However, the motion path is

mechanisms while the motor positions to be generated are discrete. One cannot generate an infinite number of motor positions. Unfortunately, the splines employed to produce the motion path do not offer convenient mechanisms to place each state with a proper interval between the states. If constant local spacing is used in an attempt to make a minimal inter-state distance, the motor positions appear as in Figure 4-1. Note that the positions do not indicate multiple positions in a band time-wise. Instead, the figure represents the positions of a single state in time piece. Consequently, when one displays these frames each made up of the given sequential motor positions the system is implicitly slowed down near the position at time t = 0, t = 2 and t = N. In the meantime, the motion speeds up and then slows down in between. The example shown here is a particular case and the relation spatial distribution of the interwoven states is even less predictable in general. As long as the unnecessary slowdown and speedup persist, as long as there is no method to put the synced motor control, one will meet with difficulty in capturing the smooth motion transition which is the original purpose of the animation.

Figure 4-2 illustrates the motor positions produced by the speed control method described in this chapter. Although it can also incorporate various offset such as acceleration and deceleration, the constant speed is applied as the example. As can be seen, the inter-motor distance is evened out to some degree, thereby reducing unwanted motion speed.

This chapter develops a method to control the motion speed, using its representation in a discrete parameter domain. Two subsequent experiment techniques, Landmark Adjustment and Averaging Adjustment, will be combined with Incremental Root Spacing to make up the motion control method.

position of time
( + )

subtantral
position

position at alpha carbon
(x axis) = 0

position at zero
(-)

Figure 9.5 A spatial distribution of six atom in time domain with overprojected quality. Between the time $(+)$ and $(-)$, $t_1 = 0$ and $t_2 = 0$, the motion is conflicted and thus disoriented upon sequential display of the atomic positions.

Figure 9-2 A spatial distribution of an orbit in lane domain when the speed control method described in this chapter is used to generate a constant speed.

## Literature Review

Association on the context of computer graphics can be clustered in two categories. In keyframe animation, the motion is determined from a series of interpolatory positions between pairs of key frames while in dynamic simulation, the motion is determined from a physical law. Given a pair of key frames, generation of arbitrary interpolatory scenes to other solid substantiating and the frames thus generated are called inbetween frames. Usually, the inbetweening is done between a pair of poises belonging to the interconnecting key frames, and the inbetween frames can be constructed from the collective interleduce poises making up a soft frame. Our concern is this chapter is with keyframe animation, particularly in the context of using in the automatic inbetweening system.

In most early keyframe animation, linear interpolation was used. For instance, in the animation system developed by Burtnyk and Wein [Burt71, Burt14], an example is given about the transformation of an alphabet to another alphabet character. If the character 'A' is to be changed into the character 'B', then both characters are decomposed into equal number of line segments so that an endpoint of the line segment in 'A' can be mapped onto another endpoint in 'B' in a one-to-one manner. Because the two endpoints, the motion path is assumed to be linear. The primary reason the linear motion path is emphasised here is that through simple means, the motion speed can be constant with constant motion or parameter t domain. Nevertheless, linear interpolation cannot produce a smooth motion path and sometimes it even

shows distortion of an object which is naturalist [Kochò]. Thus it was abandoned in favor of a more natural cubic spline interpolation.

In subsequent development, a real object was represented with a simplified form, a rock figure or a skeleton [Hoffö]. Since a skeleton can be drawn easily, the more skeleton drawn between key frames take the role of additional intermediate key frames as they first the motion speed and the motion path can clearly be controlled. Although skeleton prevention curve filling [Akaïö] is used in this scheme to smooth the motion path, the motion speed and motion path in this scheme are mostly controlled by an interactive tool and over procedure in the lowest level from which no further information will take place. Since the user is required to manipulate the interpolated frames, the burden of kinematic adjustment is left to the user in practice. Further, the polygons enclosing the stick figures still were associated with linear artifacts [Burtlö] in this method.

Another approach that requires a user specification down to the level of inbetween frames can be classed in GRASS03 [Baerö]. Since the animator provides a continued description of the path and timing of a motion by drawing motion paths called P-curves on a digitizing tablet. The shape of a P-curve defines the inbetweening trajectories of the motion path. In practice, a trail of symbols is read instead of a continuous line to depict the path, and the symbols represent where the inbetween frames should be put on the motion path. In addition, the symbols are spaced equally in time so that the dynamics are represented by the local density of the symbols. Therefore, control of the motion speed, which is represented by relative spacing of the inbetween frames, should be dependent directly on the animator's intuition.

The inverting power movement in Barzer approach [Baerö] offers a further degree of control in the P-curve method. A moving point is defined

much like the symbols in a P-curve: a curve in space and time which constrains both the trajectory and dynamics (i.e., path and speed) of a point. While the P-curve is limited on the motion of a single point, the scheme allows the specification of the motion of multiple points. For instance, a curve or an object of animation can be approximated by multiple points lying on it. If an animator specifies the motion path for each of these points, the set of key frames and the set of moving points specify a patch network of the motion sequence. Then the set of key frames in a single patch is automatically interpolated using such methods as in [Milnet, Cox74]. However, the motion speed in this method will tend to be manually adjusted.

Steketee and Badler [Ste85] addressed the problem of the motion of the temporal aspect from the spatial aspect of animation so that each can be put under separate control. The control of these aspects was made possible by the composition of two functions, the kinetic comparison and the position interpolant. The kinetic interpolant expresses the keyframe sequence as a function of time, and returns an information about the actual values of the motion parameter (i.e., joint parameter) that determines the position of the key frame. The position interpolant expresses the spatial position defined by the key frames, as a function of keyframe sequence. Therefore modification of the kinetic interpolant changes the timing of the key frames, and hence the speed and acceleration of the motion, without causing any change in the spatial position of the key frames. In other words, the motion speed was arbitrarily controlled by adjusting the time interval between the display of successive key frames, whose spatial position is fixed by the position interpolant. Although the function for frame time was invisible to the animator to modify, the positional change behaved by the modification of the kinetic time cannot be intuitively grasped with this approach.

Burnds and Haddon [Bartlett] allowed direct control of the motion speed by adjusting the sampling speed in the laser parameter domain. Their approach was aimed at providing the estimator a means of a direct manipulation of the object-space speed rather than the indirect control of speed achieved through manifestations of frame times. In this approach, the desired magnitude of the speed (i.e. distance between the estimates framed on the parameter curve $P(s)$, as a function of a motion speed parameter $s$, is defined as a quad-profile curve. On a parametric curve $P(s)$, the least parameter $s$ can be related to a curb that equally spaced values at a well produce values of $s$ that lead to the desired sampling pattern (i.e. motion speed) at the arbitrary known points for $P(s)$. As the function relating the least parameter $s$ and the speed parameter $s$ needs to be determined. If we define this function as the/test profit, then the problem of control of motion speed in their approach reduces to solving a differential equation about the least profile defined in terms of the profit curve and the spline function adopted for the interpolating.

However, there are a few problems associated with this approach. First, finding a sequence of the least parameter $s$ requires iterative evaluation of the Newton-Raphson integration [Newton], with an inconvenience in terms of the computational speed. Second, the root finding problem inherent in the boundary condition in the deferential sequence is difficult to handle in some cases, because it intrinsically has a tendency to diverge, depending on the situation [Prett]. Third, this method does not necessarily result in 'hitting the key frames.' That is, the key frames may not be chosen during execution. The effect to control the spatial distance between the (determined) frame estimates key frames from the estimator domain [named]. The main reason for

the loss of the key frames means here the cumulative errors associated with numerical evaluation of the differential equation.

The approach taken in this chapter shares some ideas with the above approach. For instance, we adopt the speed profile curve, and use a similar relation between the motion speed and the low-frame distance. However, our approach removes the evaluation of the differential equation and replaces it with simpler subroutines followed by algorithmic adjustments. The direct interpretation of incorporating the algorithmic adjustment procedures is that a generation one will fit the key frames while providing appropriate motion speed to the underscored frames. The work done in this chapter is required, once our prototype system described in Chapter 5 is based on the assumption that the key frames must be part of the animation sequence.

## Incremental Knot Spacing Method

The principal difference between Bartels and Hardtke's approach [Bar91] and our approach presented in this section lies at the method for the determination of the knot parameter $t$ that controls the motion speed. While they employed the knot parameter $t$ as a continuous medium employing Budge-Kutta integration, our subroutine replaces it with a technique called incremental knot spacing. This technique, combined with a few adjustments described in subsequent sections, plots the control of motion speed on a simple subtraction basis. This section describes a method to calculate the knot parameter $t$ that could indirect the control the motion speed in animation.

Although a rough concept of the speed profile curve can be found in [Bar91, Men04, Ste92], it tends to be more precisely defined in our approach,

d in the relative distance between the inter/manual frames as a function of frame sequence number

$$M_i = b \, l s_i^f \tag{6.1}$$

where

$M_i$ is the relative inter-frame distance (i.e., speed),

$s_i$ is the frame sequence number (i.e., $s_i = 1, 2, ...$),

$b$ is the mapping that assigns a desired speed value for each frame number.

For instance, the value of $M_i$ represents if a maximum speed is desired. If five inter/manual frames need to be produced with constant speed, the $M_i$ values for each frame may be all 1's. However, assignment of all 1's still generates a constant speed is not approach, and thus the term relative distance instead of distance is employed in the definition of the speed profile curve. Since the absolute speed at a certain frame does not have any meaning by itself, and the speed relation to other frames reflects the spatial positioning of the frame on the curve, the definition is valid without loss of generality. It should be noted that this definition differs from the definition of speed in a common sense. Loosely speaking, speed as this context is the distance traveled per unit frame, while physical defines the term as the distance per unit time.

It is important to note that the interanimate definition of the speed control parameter is is slightly different from ours. In the interanimate definition, the values of $n$ are assumed to be combination and an explicit relation at animated between the a value and the frame sequence number. In our approach, the $n$ values can be only discrete integers (i.e., $s_i = 1, 2, ...$) representing the frame sequence numbers. The reason for this discretization

of it is connected to our subsequent approach to determine the determination of the knot parameter. In addition, the speed profile curve becomes more intuitive since the profile represents the speed of the frames instead of the speed in relation to the motion parameter u. Figure 4.8 is an example of a speed profile curve. The speed values of the sequential inbetween frames numbered 1 through 8 are shown with corresponding magnitudes, $u_1$ through $u_8$. Note that when we define as a speed profile is the shortened version of the conventional continuous speed profile curve (dp). Hence the continuous profile can be approximated by the discrete profile simply by sampling the continuous speed profile curve with equidistantly spaced u values where the number of sample is equal to the number of inbetween frames.

Moreover, the speed values between the key frames can be interpolated such that the inbetween frames can be assigned interpolated speed values. There is no constraint on the reproduction scheme in selecting the proper speed. The scheme to control the speed may be a central order polynomial to enable continuous and deceleration. Figure 4.9 is an example of the use of a b-spline to interpolate the speed values several key frames. Our advantage of interpolation of speed values is that it smooths and spreads out the speed values so it does to the interpolation of speed into person. Whatever the continuous profile curve is, the discrete speed values can simply be taken by sampling. For instance, the $u_3$ figure 4.8 is the discrete speed value corresponding to the inbetween frames with sequence number 3.

Given the desired speed profile along with the interpolation scheme, the conventional approach is focused on finding a function that relates these two input parameters in the determination of the knot parameter u

$$s = s (u)$$

such that the equally-spaced values of $n$ will produce values of $s$ spaced along some desired pattern of points on the trajectory $P(s)$. To distinguish these curves from the spiral profiles we use the term knot profile as a function from the horizontal axis to the knot profiles in hyperplane of this hardware as inflation. Consider the case of constant speed in which the distance between the inbetween frames needs to be identical. In other words, we wish to maintain constant distance between the sample points as curved on $P(s)$. According to the conventional approach, the parameter $n$ is supposed to advance with constant increments. Thus the knot parameter $t$ will be evaluated for each of these $n$ values using the above equation. Finally, each evaluated knot parameter $t$ will be plugged into the curve representing $P(s)$ to generate the sample points with constant distance. Note that the knot parameter $t$ does not increase with a constant increment even though the parameter $n$, and the final distance between frames does. Therefore, the problem of control of the motor speed is reduced to finding the function that could relate the parameter $n$ to the desired sampling pattern on $P(s)$.

Since, in the conventional approach, the knot parameter $t$ is a continuous function of the motor speed parameter $n$, evaluation of $t(n)$ requires a numerical approximation where $n$ values need to be evaluated for more densely (e.g., by around a factor of 10) than the number of in-between frames, and still we would inhibit cumulative errors. The method is not costly, particularly when a local evaluation or deceleration needs to be approximated.

Instead of using the above equation to find the knot parameter $n$, Incremental Knot Spacing sets the incremental difference in the knot

parameter $x$ to be a function of the dimensioned spread profile $\psi$, with $k-1$ number of knots),

$$\Delta x = x_{i+1} \cdot x_i + f(\psi, \delta_i) \qquad (4.2)$$

The subscript $i$ designates the $i^{th}$ subvertex frame such that $x_i$ produces the $i^{th}$ subdomain frame upon insertion into the spatial interpolation curve $f(\psi)$. If a knot parameter $x_i$ is determined for one frame, the next knot parameter $x_{i+1}$ is determined utilizing a function of the current knot value $x_i$ and the spread profile value $\psi_i$. Note that by defining this way, the recurrence relationship between the spread profile variable $x$ and the spread parameter $x$ can be modeled and bases numerical integration becomes unnecessary.

The problem of the parametric representation of a curve lies in the fact that the incremental chord length is not constant with respect to its figure left; the distance represented by $\Delta x_i$ is the spatial distance between the position of the $i^{th}$ frame at $P_i$ and that of the $(i+1)^{th}$ frame at $P_{i+1}$. Both positions $P_i$ and $P_{i+1}$ are mapped from the knot parameters $x_i$ and $x_{i+1}$ onto three-dimensional space and thus the increment is the knot parameter domain. As is implied into the spatial displacement must change; the local parameter increases with spacing $\Delta x_i = x_{i+1} - x_i$, though the local parameter increments with constant steps $\delta_i = \delta$, i.e. the corresponding increment in the space, $\Delta x_i$ varies in general. The subsection that a uniform process or parameter domain yields non-uniform spatial distribution in space domain, is typical at the behavior of generic cubic splines where

$$P_i = f(x_i), \quad \delta_i = x_i$$

$$\Delta x_i + f(x_i) \cdot x_{i+1} \cdot \psi_i$$

The increments chord lengths between any two adjacent points is

$$\Delta l_i = \left| P_{i,n+1} P_i \right|$$

$$= \sqrt{\sum_{m=1}^{3} \left( P_{i+1,m} P_{i+1,m} \right)^2} \tag{4.2}$$

where

the summation takes over the $n$, $x$, $y$, $z$ components of the three-dimensional Euclidean space, and

$P_{i,k}$ is the $k$-component of the processed vector $P_i$.

Therefore, evaluation of the distance between feature vectors is a square sum of a node under polynomial about joint parameter $s_i$ whose coefficients depend on the interpolation scheme. Because of that complexity, it is difficult to predict the behavior of $\Delta l_i$ upon constant decrease in the parameter $s_i$ to $s_f$. Moreover, a relatively large increase in the parameter domain can produce a relatively small spatial displacement between frames. Therefore, the attempt to control the speed $\Delta l_i$ by adjusting $s_i$ through the use of direct expression fails.

In our method, we concentrate on the tangent information extractable from a given composition scheme. In differential geometry (Guo58), the arc length of a parametrized curve $P(t)$ with normalized parameter $t$ is given by

$$L(s) = \int_0^s \left| \frac{d P(t)}{dt} \right| dt \tag{4.3}$$

with the tangent vector

$$\left| \frac{d P(t)}{dt} \right| = \sqrt{ \left( \frac{dx(t)}{dt} \right)^2 + \left( \frac{dy(t)}{dt} \right)^2 + \left( \frac{dz(t)}{dt} \right)^2 } \tag{4.4}$$

This is true for any parametrization including the parametrization by chord length. One consequence of Equation 4.4 is that the tangent vectors can be

used as a measure of distance between adjacent admissible frames until it is the rate of the change in arc length or the chord length of the admissible frames are densely spaced. For instance, if the magnitude of the tangent vector at a certain frame on a curve is large, the position of the next frame is inclined to be farther apart spatially. If the tangent vector at point $T_i$ is denoted by $\mathbf{d}_i$,

$$\mathbf{d}_i = \frac{\partial \phi(t)}{\partial t}\bigg|_{t=t_i}$$

with

$$|\mathbf{d}_i| = \sqrt{d_{ix}^2 + d_{iy}^2 + d_{iz}^2}$$

$$d_{ix} = \frac{\partial \phi_x(t)}{\partial t}\bigg|_{t=t_i}, \quad d_{iy} = \frac{\partial \phi_y(t)}{\partial t}\bigg|_{t=t_i}, \quad d_{iz} = \frac{\partial \phi_z(t)}{\partial t}\bigg|_{t=t_i} \quad (6.2)$$

Figure 6.8 shows a tangent vector $\mathbf{d}_i$ at a point $P_i$, a linear position in space produced by stepping from the parameter $t_i$. As is the case with admissible geometry, angular point heading to a $|\mathbf{d}_i|$ value of zero should be avoided in the evaluation of $\mathbf{d}_i$. Since the zero tangent means that the frames overlap at the same position, this situation does not occur in the natural interpolation process. To further understand the importance of the tangent vector in cubic splines, the magnitude of the vectors are examined with constant increase in the parameter $t$.

Figure 6.9 shows the position of inbetween frames based on the two-dimensional key positions in Figure 6.8. The square line joins #1 key frames #3 and Key frames #5, Key frames #6 of Figure 6.8 are interpolated using an arbitrary method spline. Nine inbetween frames are generated between key frames #3 and #5. Figure 6.9 also illustrates the variation of the magnitude of tangent vector $\mathbf{d}_i$ with constant stepping of the limit parameter (i.e., $\Delta t_i = 1$).

Although the parameter $t$ increases with contours spacing, the distances between frames are unequal. In practice, the magnitude of the largest vector in a second-order polynomial of the parameter $t$ with its coefficients determined by the positions of four nearby key frames.

It is important to note that the shape of the curve itself does not define the tangential magnitude at a certain point on the curve. Although different interpolation schemes, and thereby a different mathematical representation may produce the same curve shape, the tangential magnitude at a point of the curve varies depending on the interpolation scheme. In other words, the magnitude of a tangent is not determined from the shape of a curve. It is not an inherent property of the curve. In fact, the nature of the tangential property of the parametric representation of a curve can be distinguished from the non-parametric curve. For instance, the tangent of two-dimensional non-parametric curve can be represented as $dy/dx$ while $x$, $y$ being Cartesian coordinates. The magnitude of the tangent can neatly be visualized such that it is large with steep slope and vice versa. This is not true for parametric curves, since the tangential magnitude is a function of the underlying implicit parameter $t$. Hence the visual slope, $dy/dx$ on a parametric curve, does not necessarily coincide with the magnitude of a tangent.

However, it is possible to relate this non-uniform tangent magnitude to the inter-frame distances to achieve control of the motion speed. Since our approach is based on the determination of the speed profile (Equation 6.1), the parameter (Equation 6.2), and the speed limit (Equation 6.3), the continuous relationship in Equation 6.4 can be easily converted into a discrete form.

$$\dot{s}_k = |\dot{A}_k| \cdot \Delta t_k \tag{6.6}$$

$$L = \lim_{\Delta m \to \infty} \left[ \sum_{i=1}^{m} \Delta L_i \right]$$

$$+ \sum_{i=1}^{n} \Delta L_i \tag{57}$$

where

$L_i$ is the total arc length between key frames,

$\theta_k$ is the tangent at the $k^{th}$ inbetween frame whose position is determined by the interpolation function $P(t)$ at $k = t_k$,

$\Delta L_i$ is the distance from $\overline{PR}$ to $L=\overline{PR}$ arcbetween frames.

The direct measurement of Equation 4 is that the chord length between frames is directly proportional to the magnitude of the tangent and the incremental parameter spacing. For anatomy, in Figure 4-5, we note that the chord length (i.e., inter-frame distance) was proportional to the magnitude of the tangent each constant incremental parameter spacing. This means that, to maintain a constant speed, the parameter spacing should be reduced if the tangent value is relatively large.

Through use of Equation 4-7, the total arc length between key frames can be approximated by the sum of the piecewise line segment $AL_i$ and this continuous approximation the total arc length as the number of inbetween frames increases. Figure 4-6 illustrates an approximation of the total arc length between key frame $k$ and $P$ by three piecewise line segments. As the inbetween frame 40 (LF 40), the incremental arc length $\Delta L_4$ can be roughly approximated by $\frac{1}{2}\theta_k |\Delta t_4$. The direction of the tangent $\theta_k$ approaches more closely the vector $(LF 40-LF 42)$, as the number of inbetween frames increases. As a matter of fact, the proximity between the chord length and arc length of the segment, what one generates as the speed is an animation.

sequence to the chord length between frames and not the arc length. What matters, however, is how to control the incremental chord length to achieve proper speed.

In order for the incremental chord length, $\Delta l_i$, to incorporate speed information, three constraints should be applied as follows:

1. It should be proportional to the total arc length.

2. It should be proportional to the desired/end speed value assigned to that frame.

3. The summation of the individual incremental arc length should be close to the total arc length as the number of inbetween frames increase.

To satisfy these constraints, the expression for the incremental chord length is

$$\Delta l_i = 1 \cdot \frac{z_i}{\displaystyle\sum_{i=1}^{n} z_i} \tag{6.10}$$

where

$1$, is the total arc length,

$z_i$ is the discretized speed value at the $i^{th}$ frame,

$n$ is the number of inbetween frames including the beginning key frame

That is, the discretized speed value is normalized by the sum of individual speed values, as first Equation 4.9 is valid upon reasoning both sides of Equation 4.6 posit is chosen.

Combining Equation 6.8 with Equation 6.6, we get the incremental knot spacing at the $i^{th}$ frame as,

$$\Delta t_i = t_r + \frac{t_s}{|\Delta i|} \frac{\Delta i}{\sum_i \Delta_i}$$

(4.9)

with the input

$$\Phi_i = \frac{dP(t_i)}{dt \cdot |\Delta i_i|}$$

where

$t_r$ is the real arc length between key frames.

$P(t)$ is an interpolation function.

The incremental knot spacing of the $i$th frame is set to be inversely proportional to the tangents of the curve in that frame. The magnitude and direction of the tangents vary depending upon the interpolation scheme and the position of the key frames to be interpolated. However, in view of the fact that the pixel tangential magnitude occupies inter-frame distance and thus the speed, the inverse relationship between the knot spacing and the tangent magnitude cannot connect speed. This happens because, if the incremental knot spacing is diminished, the inter-frame distance tends to decrease, as can be verified by Equation 4.6. In contrast, the incremental knot spacing is set to be proportional to the discretized speed value of that frame. This setting is quite natural since the desired speed should be applied proportionally to the inter-frame distance.

Consequently, control of speed can be accomplished by adjusting the incremental knot spacing according to the desired speed profile curve. In the beginning, the knot value of 0 is set to zero to produce the first in-between frame which is practically a key frame. Then Equation 4.9 dictates the incremental knot spacing to be applied to the current $n$ value. The summation of $\Delta_i$ and the next knot parameter value corresponding to the

several unknown items, and so on. Therefore, for any two adjacent items, the distance from the first item to the second is determined from the final value of the first item using incremental item spacing.

In summary, the section has described the incremental item spacing method for speed control in simulation. Because of the difficulties involved in evaluating the displacement between frames, an approach to control speed by exploiting the largest magnitude a personoid found on the continuous representation of arc length, we have derived a discrete-domain counterpart as an approximation tool. Furthermore, the speed profile curve has been incorporated into the item spacing expression.

## Locational Adjustment

One of the problems associated with Equation 4.6 and Equation 4.9 is that it is valid only for sufficiently small values of the incremental item spacing, as: For large values of $\Delta s$, the distance between two adjacent achievement frames and between and the tangent may not properly approximate the curve behavior between the frames. This is true especially when the number of in-between frames decreases so that the total arc length is approximated by only a few line segments. The problem is undercompiling in which the each sample represents the position of in-between frames.

Figure 4.6 illustrates an undercompiling error when three in-between frames are used to interpolate a mature path P62. At in-between frame #1 (I.E. #2), Equation 4.6 was used to evaluate the distance item that locates in the next frame (I.F. #2) to extend the motion speed. Thus from location 4.5, the incremental tangential magnitude was calculated by differentiating the curve P62. There two constraints determine the incremental item spacing in

so that the knot value $v_i$ corresponding to $\widehat{G}F$ RD can be recovered. However, upon insertion of the knot value into curve $F_{34}$, the point on the curve corresponding to the knot value may not be in the desired position i.e., $\widehat{LF}$ RD. The point may be before or after the $\widehat{G}F$ RD and there is no way to predict the position of the $F_{34}$ on the curve.

This error is a result of the assumption that the tangential magnitude is constant between frames. In the above example, the magnitude and direction of the $\widehat{d}_0$ were assumed to be constant throughout the interval [$\widehat{LF}$ $R_0$, $\widehat{LF}$ $R_1$] and our method predicts the position of the inbetween frame $\widehat{F}_0$ of $\widehat{LF}$ $R_1$ which we no longer require on the curve $F_{34}$ if a tangent vector of each inbetween frame is used, and then if the interval distance is small, the tangents can properly represent every inbetweening tangent obtained from the curve point in the desired interval. Since the general cubic spline is based on polynomials which are relatively differentiable at all points on the curve, and since the tangential magnitude can be represented by the square root of the second index polynomial of knot parameter $t$, it is reasonable to assume that the tangential magnitude is constant within a sufficiently small value of $\Delta t$. However, this is not true if the inbetween frames are sparsely spaced. The larger the interval between frames, the greater the chance that the curve exhibits a change in magnitude.

Figure 4-16 shows a schematic diagram of Lambertian adjustment. The knot value corresponding to the inbetween frame $F_i$ as evaluated in Figure 4-14 is inserted into the interpolatory curve $F_{34}$ so that the positional vector $\mathbf{r}_i$ a comdate of the inbetween frame $F_i$ is generated.

$$\mathbf{r}_i = F(t_i) \Big|_{t_i = v_i}$$

Note that the trial position $P$ nor past the desired position $P_d$. This situation may arise when the tangential magnitude is changed to position $P_i$ is underestimated. That is, the actual tangent value has returned to the desired $\{B_2 P_2\}$, and the tangent at $P_2$ has failed to represent all the tangents at that interval. Due to an underestimation of the tangential magnitude, the incremental laser spacing $\Delta x_{i}$ was overestimated and thus the trial position ran past the position of the final inbetween frame $P_d$.

Instead of using the trial position $P_3$ as the position of the frame, forehand adjustment uses the trial position as a basis to generate the proper frame position. In Figure 4-8, if we denote the distance laser position $P_i$ to trial position $P$ as $\delta t_i$, the general expression for the relation between them is

$$\delta t_i = |P_{i+1} P_{i}|$$

where

$\delta t_i$ is the distance from the $i$th frame to the $(i+1)$th trial frame,

$P_{i+1}$ is the position of the $(i+1)$th trial frame,

$P_i$ is the position of the $i$th $i$th frame.

One way to express this information about the distance from the current position to the next trial position is by scaling. Let the laser spacing originally used to estimate the trial position be $\delta t_i$, and the new laser spacing to be evaluated be $\delta a_i$. Then we represent the new laser spacing such that Equation 4-6 is still valid.

$$\delta t_i \quad \delta a_i = \delta_i$$
$$\delta t_i \quad \delta t_i$$

where

$\delta_i$ is the tangential magnitude at the $i$th frame.

With respect to $\Delta\sigma_i$,

$$\Delta\sigma_i = \Delta\sigma_j \frac{\Delta L_i}{\Delta L_j}$$ (6.18)

If we regard the ratio of the original chord length to the real chord length as a scale factor, then the new local spacing is the real knot spacing value multiplied by the scale factor. If the original chord length specified by the speed profile curve is the same as the real chord length, then the scale factor is unity and the required knot spacing value is maintained. If, on the other hand, the real chord length is greater than the original chord length, the new knot spacing value is scaled down, reflecting the difference in chord length.

Equation 6.10 can be mentioned with Equation 6.9 for an evaluation of knot parameter value that could control the speed of animation more accurately. From the position of a current frame, the next frame position can be evaluated in the following way:

1. Equation 6.9 is used to get the incremental knot spacing in the current frame.

2. Incremental knot spacing is added to current knot value to produce a real knot value

3. The real knot value is inserted into the expression for the interpolatory curve to yield the real chord length.

4. The real chord length is incorporated into Equation 6.18, giving rise to a new incremental knot spacing.

5. The new incremental knot spacing will be added to the current knot value and the result will be inserted into the interpolatory curve to produce the position of the next frame

As the number of unknown feature sentences, the errors at the individual intervals between frames can also be distributed by this approach. Although the errors due to constant tangent representation at a certain interval can be greatly reduced with this approach, there still remains the problem of errors between the junction calculated from one approach and the junction dictated by the profile curve. While this section has dealt with reducing the positional errors between two adjacent frames, the next section will treat the cumulative errors associated with the total references frames between a pair of key frames.

## Averaging Adjustment

Using Equation 8.9, the spatial distance between adjacent frames and thus the speed between frames can be calculated based on the speed profile curve. In particular, if the motion pitch between a pair of key frames is foreshortened by a limited number of frames, application of a uniform tangent value between adjacent frames leads to the error in the spatial distance. This error is defined to be the distance gap between the references distance calculated by the incremental knot spacing and the one specified by a speed profile curve. Although Lookahead Adjustment could reduce the error between any two adjacent frames, the remaining errors will accumulate throughout all the interconnected frames. A direct consequence of this cumulative error is that the key frame might be moved and it may not be displayed as part of an animation sequence. Since our approach is based on the principle that the key frame must be displayed during animation, the problem must be handled properly. This means another fine problem by redistributing the cumulative errors into individual frame distances.

Before delving into further discussion, it is important to re-examine the relation between the bonus and the normalized lost parameter $x$. Given a pair of key frames to be interlocuted, the position of the first key frame is assigned to be mapped from the parameter value of zero while that of the second key frame is assigned to be mapped from the parameter value of one. The assignment are made possible by applying these constraints into the expression for the spline representing the motion path. For instance, if three adjacent key frames $P_i F_1, F_2, R_3$ together with a pair of splines $P_i(s), F_2(s)$ are given, the later parameter $s$ will run from zero to one till in the region $[F_1, F_2]$ based on the function $F_1(s)$. Once again, the value of the later parameter $s$, which varies between zero and one, will be mapped into $F_2(s)$ to produce the interweaved later position. Although $R_2$ belongs to both splines, it is normally treated as part of the second spline. Therefore, interlocuting a pair of key frames starts with a value of zero, which is the first key frame, and ends with designating the last interweaved frame position right before the second key frame.

Since Incremental Keys Spacing is based on the chordon knot spacing $dr$, and since the knot parameter $r$ is assumed to be normalized, the root of $dr$ should be exactly zero. Only if this state is identical to zero, can the second key frame be properly interlocuted in reference to the last interweaved frame, although the actual generation of the key frame is performed through the evaluation of the spline or the subsequent interval. This is true since the second key frame is assigned to be mapped from the knot parameter value of one in the normalized parameter representation of a cubic spline. The reason why Increasing Adjustment is needed is because the root of the chordon knot parameter spacing $dr$ tends to deviate from one as the number of interween frames is reduced.

Figure 4-11 illustrates the Averaging Adjustment process. The positional vectors of the Nave inbetween frames G# #1 through I.F #20, excluding the first key frame A.I.F #3, are represented by $P_i$ with $i = 1, 3, N$.

The frame positions calculated by Incremental Keys Spacing followed by Lookahead Adjustment are shown in $P_i$, while the frame positions resulting from Averaging Adjustment are denoted by the positional vectors $P_i$. Consider the position $P_j$ of the inbetween frame I.F #3 which is the last inbetween frame preceded by the spline Plot. The key frame 0 is treated as the first inbetween frame of the subsequent interval generated with the least parameter $s$ value of zero. Therefore, Incremental Keys Spacing method will strive to evaluate the next frame position at point $P_j$. If, however, the method continues to evaluate the next frame position $P_j$, the calculated position over to the key frame 0. Since Incremental Keys Spacing method in this case has divided the range of the normalized parameter $s$ (i.e., $0 \le s \le 1$) into three subintervals $\Delta s_1$, $\Delta s_2$ and $\Delta s_3$, the least parameter $s_j$ is obtained by adding $\Delta s_1$ to $\Delta s_2$ should be taken so that the last frame 0 is assumed to be mapped from a value of one, the least parameter $s_j$ should become one to fulfill its constraint.

Unfortunately, the scan of the incremental keys values and thus the value of the least parameter $s_j$ generally does not become one. Due to the accumulation of errors between adjacent frames, the should-be one least parameter $s_j$ may extend well below one. The extra of the least parameter $s_j$ does not become one is filled up and the distance between frames as inbetween and the least parameter $s_j$ becomes the least parameter $s_j$ value. The most of this section is addressed in the next section.

The difference between frames in the frame positions as assigned proper least parameter $s$ values. The errors in parameter space are propagated and accumulated, since the value of $s_j$ is inherited from the $s_k$, and so on. The cumulative errors in the parameter space are represented by the parameter value of key frame 0

Figure 4.11 shows the case when the parameter $z_i$ is much greater than one. It also maps this value into spatial position via a spline Plot, for corresponding position $P_a$ is far away from the desired position of the key frame B. In a similar sense, the spline Plot with normalized parameter $z_i$ is stretched outside the interval [0, 1], but a value still exists since it is a polynomial. However, we do not need to evaluate the position at our approach.

Averaging Adjustment and Indicates this cumulative error into individual inter-frame distances between adjacent frames. That is, if the should-be-the parameter value $z_i$ exceeds one, the preceding increment $s_1, s_2, \ldots$ will be reduced. Similarly the parameters will be scaled up if the parameter corresponding to the key frame B falls short of one. In doing so, we scale the incremental knot values to detect proportion their magnitudes to the sum of adjusted incremental knot values becomes one.

$$\Delta s_n = \frac{\Delta s_n}{\displaystyle\sum_{i=1}^{n} \Delta s_i}$$ (4.11)

where

$n$ is the number of in-between knots.

Geometrically, scaling the knot parameters this way will shift the position of the key frame to the left or right of the previous frame positions. Figure 4.11 is an example of when the frame positions are shifted left, reducing the reduction of the incremental knot value. Note that the correct position of key frame B now can be produced with a value of one, although the reduction of this key frame is still left to the interworking of the subsequent segment.

It is by this adjustment of the incremental knot values such that they run up to one, that the key frames can properly incorporated as part of an

animation sequence. For example, in a sequence of three key frames, each key frame can be seen by the animator at a value of zero in the corresponding spline expressions. Nevertheless, the location of the last inbetween frame is apt to be discerned without foretuning. Adjustment in Figure 6-13, the last inbetween frame position $F_3$ appears key frame 3 to lie in position $F_4$, the first of the animation sequence by inserting into it into the next interval, the resulting distance between the last inbetween frame and key frame 8 is far from the expected speed as sketched before. By relaxing and computing and the cumulative error at the portion $F_4$, the transition from the last inbetween frame and the next key frame is stretched out.

Because of the binding of the transition, our approach can span across multiple key frames. While the interbetweening in Hax90 treated entire keyframe sequences as a single pair of key frames, starting with the first key frame and ending with the last key frame, our approach splits the entire keyframe sequence into subjacent pairs of key frames and provides proper transitions between keyframe intervals. The entire advantage of such a splitting is that each pair of key frames can be filled with inbetween frames with independent speed control without affecting the speed of nearby key frames. Such control between a pair of key frames with discretized speed values is interweaved while varying the normalized time parameter between zero and one, and the burden of interbetweening entire key frames with normal key frames can be avoided. These assets will be explored further in conjunction with the subsequent section about performance evaluation.

Although Averaging Adjustment does enable smooth transition to the next key frame, there is a price to be paid for such convenience. We cannot foretell the sum of the incremental time values a priori. The denominator of

Equation 4.21 is now used and the Incremental Rises Spacing method completes evaluation of the entire series of reference positions. Therefore, it is necessary to let the Incremental Rises Spacing run in the first part of the program and the Averaging Adjustment take place in the second part. The next section will show how to implement this procedure in detail.



Figure 4-9 A speed profile curve that represents speed values (i.e., interframe distance), the h-axis number is run idem electron values only

Figure 4.4 A B-spline curve used as an interpolation scheme for the speed. The speed $q_1$ through $q_5$ represents the speed values at corresponding key frames $K_1$ through $K_5$. Note that $K_3$ does not denote spatial position. In addition, the interpolated speed values at key frames do not coincide with the specified value because of the approximating character of the B-spline.

Figure 6-5 Relation between $\Delta z_i$, $\Delta z_0$ and $d_i$, $P_1$ and $P_0$ are positions of a pair of key frames in three-dimensional space; $P_i$ and $P_{0i}$ are the positions of two consecutive key frames, $d_i$ is the distance between $P_i$ and $P_{0i}$; the cubic spline interpolation with corresponding knot parameters is a smooth line. The vector $A_i$ is the tangential direction at $i^{th}$ frame positioned at $P_i$ on the curve $P(s)$.

Figure 8-6 A template location of the key frame K.F. 4 denotes the position of the key frame 4 in two-dimensional space. For illustration purpose the regions between K.F. #1 and #3, and K.F. #3 and #4 will be shown in subsequent examples.

Figure 6-7 Variation in the magnitude of the tangent vector $|\dot{x}_i|$. The black dots represent key frames while the grey dots represent the positions of the interpolated (intervowen) frames. The introduced parameter $s$ runs from zero to one with control correction ($t_i = \delta, t_i' = 1$), while the frames run from the key frame 42 to 40. The magnitudes of tangent vectors are normalized such that the highest value $|\dot{x}_i| = 15$ has the value at 1.

Figure 8-8. Approximations of the total arc length by the sum of incremental chord length $\Delta l_k$. The spline P(s) interpolates key frames A and B. key frame A is numbered as interstream frame #1 ($L\# \#1$). The total arc length from key frame A to B can be approximated by the sum of $\Delta l_k, (k = 1 - 7)$.

Figure 6-9. Approximations error due to undersampling. Keys: valued $\mathbf{q}_k = \mathbf{\mathcal{Q}}$ and $(n + 2)$ map into key frames A and B, respectively. These laborious frames including key frame A were generated to interpolate key frame A and B. The laborious frame LF #2 was mapped from the least parameter $r_2$. The distance from LF #1 to LF #3 was determined to be $\Delta r_3$ by a speed profile curve. The tangential magnitude at LF #2, $\mathbf{q}_2$ is calculated from the differentiation of the curve $P_2(s)$ with respect to $s$. Based on these constraints, $\Delta r_3$ was calculated using Equation 6-3 and added to $r_2$ to produce $r_3$. However, the frame position mapped from $r_3$ is not necessarily located at LF #3 lying on the curve.

Figure 4-11 The Levinsohn Adjustment. Knot values $k_2 = 59$ and $k_4 = 79$ map into key frames A and B, respectively. The positional vectors $P_2, P_4$ represent the position of inbetween frames 69 and 83. The positional vector $P_3$ the point on the curve P63 produced by the outnominal knot spacing 63 ($k_3$ part 63), could be further adjusted to be moved knot point $P_3$ by a Levinsohn Adjustment.

Figure 4-11: The Averaging Adjustment. Original frame positions $P_i$ and the linear positions resulting from Averaging Adjustment $\tilde{P}_i$ are produced from the corresponding knot parameter values $u_i$ and $\tilde{u}_i$. The spline $F(u)$ interpolates the feature path between key frames A and B. Key frame A is treated as the first subframe frame so mapped from a knot parameter value of zero.

## Pseudo Code Segmentation

Incremental Knot Spacing, Lookahead Adjustment, and Averaging
Adjustment together make up our approach to controlling the motion speed
or animation. Although each of them has been expressed explicitly using
segments, we need to represent the methods from the point of view of
programming to further clarify the approach.

It can be said that Lookahead Adjustment first treats the constrained
knot values locally while Averaging Adjustment reconstrues the incremented
knot values in a global way. At a certain position of an unknown frame,
Lookahead Adjustment predicts the next frame position thus redefining one
knot deviation from the distance specified by the speed profile curve.
Therefore, it applies to the distances between two adjacent frames. On the
other hand, Averaging Adjustment applies to the positions of the entire
series of inbetween frames.

Although Averaging Adjustment must be run in the second pass of a
program, the previous knot values and the incremented knot values can
be directly entered from the first pass of the program. Therefore, most of the
procedures contained with the first pass, such as the evaluation of speed,
chord length and tangents, can be avoided to speed up the second pass. The
following Pseudo Code presents code to made up of four parts: input
specifications, output specifications, the first pass, and the second pass.
Incremental Knot Spacing and Lookahead Adjustment are incorporated into
the first pass. In contrast, the Lookahead Adjustment shows a possible
iteration to further reduce speed errors. Lookahead Adjustment is

Incorporated into the second pass of the program titled as 'Route Controlled by The Speed Profile."

Program RouteControlledByTheSpeedProfile;

BEGIN

Input:

PassList – spatial position of four sample data points

TypeTable – a spline type to determine shape

TypeSpeed – a method to assign a speed profile
          here is a function at a SpeedList with a spline type to
          interpolate the speed

WaveFrames – number of inbetween frames

NumIterations – number of iterations for Lookahead Adjustment

Output:

The spatial position of the inbetween frames including the
beginning key-frame: Spline[ ]

Next pass:

```
FOR i=1 to (NumFrames-1)
DO BEGIN
      KnotList[i] = i/NumFrames
END   /* initialize list is parameter with equal knot spacing */
TotalChordLength = KnotChord(KnotList,TypeSpline)
      /* sum chord length with creation de_*/
```

TotalSpeedQuantity = SortSpeed·Open(EstClrTypeSpeed,SizeList)

/* case A, mapped from best ticket queued from contest A[ ... ] */

ScanK = 0

/* the variable to be passed into the second pass is initialized to zero */

FOR i=1 TO NumFrames

DO BEGIN

/* Incremental Reset Opening method follows*/

Si = ResetList[i]

CurrentLocation = GetLocation(TypeInflow,NoiseList,Si)

/* Find the frame position matching Si */

Distinction = GetExecutionType(InflowTicketList,Si)

/* Find the magnitude/ magnitude [ ... ] at current position */

Speed = GetSpeedType(Speed,SpeedList)

/* Evaluate the theoretical speed value by a function or direct
assignment  */

DoEach = (TotalCheckLength/Distinction) *
(Speed/TotalSpeedQuantity)

/* incremental knot opening is determined  */

FOR j=1 TO NumEntrance

/* the comparisons between the total check length and the
specified speed profile can be iterated */

DO BEGIN

TimeSik = Si·DoEach

/* Get the knot value of the next frame */

TimeLocation = GetLocation(
TypeInflow,NoiseList,TimeSik)

/* Find knotMaked position mapped from $t_{Sik}$*/

```
        LocationDifference = AbstractValue
                           (TextLocation−CurrentLocation)

        /* Extracts the error in the speed */

        Delta% = DeltaK *
                  (Elefix%*Speed%/LocationDifference)

        /* Lookahead adjustment is complete */

  END

  SilentE% = 8
  /* Save for second pass */

  DeltaSilentE% = Delta%
  /* Save for second pass */

  Iterm% = Iterm%+E%
  /* Save for second pass */

  ScanLook%% = ScvDelta%%
  /* Adjust next learn location */

END


Second Pass:

  Iterm% = Iterm%
  /* Assumed ∑ dx_i from first pass */
                  i=1

  Splineoc% = GetLocatorType(Spline,FromList,S%)
  /* First print of the format */

  FOR i=1 TO NumFrame
```

```
DO BEGIN
    Si := Sil.get()
    /* Recovers si from first join */
    Delta := DeltaM.get()
    /* Recovers dm from first join */
    Delta4 := Delta3/Sum3i
    /* Antstaging adjustment */
    Si := Si+Delta3
    /* Can we is ? */
    Spline(x) := GetLocation(TypeSplines,Found.at(Si)
    /* Locate the frame positions with the new knot values */
END
```

### Brief Analysis of the Algorithm

This section illustrates visually how incremental Knot Spacing and the subsequent adjustments enhance the control of inter-frame distances. The control of the inter-frame distances can equally be said to be control of speed, since speed in our approach is defined as the distance between frames. To facilitate the understanding of positional variations of the frames under our approach, the adjustments are tested against the implementation of the constant speed between adjacent frames.

All the figures shown in this section are based on the two-dimensional template data-position (i.e., keyframes positions) in figure 4-6. As a result, the frame positions appearing in the subsequent figures denote the positions in the two-dimensional Cartesian coordinate system. Two-dimensional data is

used only for the purpose of convenience for illustration, and the adjustment techniques are equally applicable to the finer-dimensional spatial data. The positions of the projected inference frames are shown only in the region between the three inner data points of Figure 6-2. Throughout the figures in this section, the gray dots represent the position of reference frames and the black dots represent those inference frames that also belong to key frames. The positional vector $P_i$ represents the location of the frames.

Before showing the spatial adjustment techniques, a comparison is made between the performance of the iterational method spline and the true base spline discussed in Chapter 3. The purpose of the comparison lies in the fact that the speed and thus the error-free distance is determined from the representation of the spline interpolating the motion path rather than the shape of the spline curve. The same curve can yield a different pattern of the inference frames depending on the mathematical representation of the curve.

Figure 6-2 and Figure 6-2b illustrate the location of the inference frames produced by motion knot spacing. Ten inference frames, including the first key frame at $P_1$, are generated at 1.0 in the motion path from the key frame to the next key frame at $P_{11}$. Similarly, the inference frames between the key frames at $P_{11}$ and $P_{26}$ were generated while varying the knot parameter $s$ from zero to 9 with constant incremental knot spacing. In both figures illustrate a highly non-constant motion path with less path speed for those inference frames between key frame $P_1$, although the same constant knot spacing was used to generate the inference frames in the two identical motion paths, the spatial pattern of the distribution of the inference frames is quite different. The measured cardinal spline SImAHI in Figure 6-2b showed a relatively dense frame distribution near the key frames at $P_1$, $P_{11}$ and $P_{26}$,

whds it shows opens distribution near the center of the adjacent ray frames. As a result, the motion appears to be decelerating near the key frames, when compared with the motion near the center. In contrast, the line form spline at Figure 4-10 exhibits almost uniform distribution near the key frames. The constant knot spacing has led directly to a constant spacing between the frames, thereby producing constant motion speed.

The different behaviors can be verified using the relation between the knot parameter $s$ and the spatial position of the Interna P0il mathematically. The interated cubic cardinal spline is a modification of conventional cardinal spline and is expressed as,

$$P(s) = [s^3 \; s^2 \; s \; 1] \; T \; [P_1 \; P_2 \; P_3 \; P_4]^T$$

with the tension matrix

$$T = \begin{bmatrix} -a & 2-a & a-2 & a \\ 2a & a-3 & 3-2a & -a \\ -a & 0 & a & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}$$

where

$a$ is tension control parameter,

P(s) is the curve interpolating the input $[P_2, P_3]$

with the knot parameter $s$ varies between zero and one,

so that $P(s) |_{s=s_i}$ is the interned frame position corresponding to $s = s_i$

The tangent vector of the spline is,

$$\frac{dP(s)}{ds} = [3s^2 \; 2s \; 1 \; 0] \; T \; [P_1 \; P_2 \; P_3 \; P_4]^T$$

Let us factor on the first two sentences at the right hand side of the above equation, the matrix $[3s^2 \; 2s \; 1 \; 0]$ and $T$. Even if the tension parameter $a$ is set

to one the high tension, the elevation in the first row of tension vector T do not vanish. Therefore, the term of $T$ will involved in the expression for the derivative of the curve $P_5()$ upon multiplication of the first two matrices. That is, the magnitude of the tangent at some point on the curve $P_5()$ becomes a second order function of $s$. In terms of the distance between frames, the inter-frame distance per unit $s$ is a second order polynomial of $s$.

In contrast, taking the tangent of the line lines spline (which does not necessarily belong to the $G^2$ class), we obtain

$$\frac{dP_5()}{ds} = [9s^2 \ 3s \ 1 \ 0] A [P_1 \ P_2 \ P_3 \ P_4]^T$$

with

$$A = \begin{bmatrix} -u+\frac{1}{t} & -1+\frac{1}{t-1} & v+\frac{1}{t-1} & -v+\frac{1}{t} \\ 3u+\frac{1}{t} & \frac{2}{t-1} & -2u+2\frac{1}{t-1} & v+\frac{1}{t} \\ -u+\frac{1}{t} & -1\frac{1}{t} & u & \frac{1}{t} \\ u & \frac{1}{t} & & \end{bmatrix}$$

where

$t$, $s$ are the endpoint knot parameters corresponding to point $P_2$ and $P_3$ respectively.

Each of the elements of first two rows of matrix A becomes close to zero if $s$ equals zero and $v$ approaches zero. The third row of matrix A combines with the $P$ matrix to produce the constant tangent vector ($P_3-P_2$). Therefore, the tangent vector is,

$$\frac{dP_5()}{ds} = P_3-P_2$$

Since the tangent is defined as the displacement in $P(s)$ unit axis $s$, the displacement between the positions $P(0)$ and $P(1)$ is the same as the magnitude of the vector $\vec{P}_1\vec{P}_0$. $P(s)$ is linear. Since the magnitude of the tangent is constant regardless of the unit value $s$, and since the curve is exactly linear, the inter-frame distance is always uniform with the highly-textured free form spline. Therefore, it can be said that the spatial distribution pattern of the frames in a interpolatory curve depends on mathematical representation rather than the shape of the interpolatory curve.

However, if the $n$ parameters change from the values coming the high interest, the diameter is that two rows of the matrix $A$ will be certain values other than zero. Consequently the non-zero rows will cross the $A^*$ zero and the $n$ rows to be part of the expression for the tangent. The rows holds true for the nonzero cardinal spline. Thus constant knot spacing in the free parameter domain no longer return constant spacing in the space domain. A key product mesured continuation and deceleration due to the non-uniform distribution of the frames on the motion path, thus mentaining the use of the speed control techniques developed earlier.

Figure 6-12 through 6-15 show a gradual enhancement in controlling the motion speed when a constant inter-frame distance is desired. An ordinary cardinal spline was used to generate the motion path of the figure. Figure 6-14 shows the refinement frame positions when no speed control strategy is used. The knot parameter $s$ domain is simply subdivided into equal knot spacing from the beginning frame 0 (i.e., first frame), 41, 42, 44 and 45 are mapped from $s = 0, s = .1, s = .2, s = .3, s = .4, s = .5$ each corresponding positioned vectors $\vec{P}_0, \vec{P}_1, \vec{P}_2$ and $\vec{P}_n$. Notice that the inter-frame distance varies as a non-uniform manner as a function of the parameter value $s$ since the tangential magnitude is a monoduvariate polynomial of the parameter $s$.

The positional vectors $P_1$ through $P_{18}$ are generated similarly and are shown here against the positional symmetry with respect to the frames $F_1$ through $F_9$. They will not be stirred in the table following the figures since the properties in the interval $[P_1, P_{18}]$ result evenly with those of the interval $[P_1, P_9]$.

Table 4-1 shows relevant measures of the orderstoring in Figure 4-18 as follows:

1. The column labeled "Frame Number" indicates the sequential frame position.

2. The column labeled "Tangential Magnitude" indicates the tangential magnitude $|\,d\,|$ on a given frame position is defined in Equation 4.5. The numbers are represented relative to the distance from $P_1$ to $P_9$ which is set to 297.

3. The column labeled "Raw Error Value" indicates the incremental limit value, $m_i$, used to prevent the next billerances frame position.

4. The column labeled "Measured Distance" indicates the raw incremental distance $M_{i,j}$ as measured from the frame position in the preceding figures using Equation 4.3.

As can be seen in the "Tangential Magnitude" column, the speed per arc-limit value $(s = 1)$ is narrowing from $F_1$ toward the center of $F_5$ and $F_9$ and from that arc-lengthes until the position $F_9$. The "Measured Distance" column, representing the measured lima-frame distance in Figure 4-18, roughly reflects the trend as the tangential magnitude, since the incremental measured distance are constant. The constant frame quantity is the parameter distance evidence sees constant spaced between between frames with a constant offset and a standard deviation of 1b, where the distance is measured relative to the length between $P_1$ and $P_9$, which is set to 297.12 corresponding olorerten

fringes are displayed sequentially, the speed picks up toward the middle of the interval $[P_1, P_2]$ and slows down toward the keyframe position $P_2$. Although the figure shows symmetrical speed control near the middle of $P_1$ and $P_2$, there is an angular pattern to the acceleration and deceleration of the speed in general, as can be deduced (see Equations 4.5 Depending upon the interpolation method and the horizon of the key frames, the separable magnitude varies unpredictably and no general rule exists governing the behavior of the sequential magnitude that could facilitate control of the motion speed.

Figure 4.15 shows the inbetween frame positions upon application of horizontal knot-spacing method. The total arc length $L$ used in Equation 4.9 is assumed to be the same as in the ordinary method; either, so that the mean value of the Movement Distance column of Table 4.1 in the inter-frame distance studied in this case. A more accurate value of the arc length can be computed by subdividing the arc with even line segments. However, the arc length thus computed does not necessarily coincide with the sum of the inter-frame distances, especially when fewer inbetween frames are required, and it is reasonable to exactly estimate the total arc length as the sum of the inter-frame distance with a given number of inbetween frames along the constant incremental knot value. The arc length of inbetween frames in a figure 4-16 are generated, but the method diversifies the incremental knot values rather than using constant symmetrical values to produce uniform inter-frame spacing. The speed is controlled by manipulating the incremental knot value $\Delta v_i$ and the knot value $v_{i+1}$ is obtained from $(v_i + \Delta v_i)$, corresponding to the $(i+1)$th frame. In Table 4.2, two new columns are defined as follows:

1. The column labeled 'Calculated Distance' indicates the distance

inter frame distance $\Delta L_c$ incorporating the speed profile using Equation 6.6 and 6.5.

The numbers are shown relative to the distance from $P_1$ to $P_2$, which is set to 100

3. The column labeled "Distance Error" indicates the error between the derived inter-frame distance and the resulting real inter-frame distance

To maintain a constant distance between frames, the elements of the column "Calculated Distance" are all set to a constant value 40. This setting is possible by adjusting the incremented knot values in association with the corresponding tangential magnitudes, since the calculated distance is the multiplication of the two columns. For instance, because of the small tangential magnitude in the first frame position $P_1$, incremental knot value is impulsed to be relatively large at the position. Despite the adjustment, one can notice the relatively large displacement of $P_1$ from $P_2$ compared with other inter-frame distances, as tabulated in the Measured Distance column corresponding to position $P_1$. This happens because the interval ($P_1, P_2$) is represented by a single constant tangential magnitude, although the actual tangential magnitude was increasing in that interval. As a result, the increased actual tangential magnitude was interfered with the relatively large incremental knot value to yield a large distance gap. In particular, the change of the tangential magnitude (i.e., the second derivative of the interpolation curve $P(t)$) is very rapid near the frame position $P_1$, as can be seen in the Tangential Magnitude column of Table 6.4. Incremented Knot Spacing could not catch up with the rapid change in tangential magnitude with so few frames.

The sum of the incremental heat values for the reaction is found exceeds one, the maximum value of the normalized heat parameter. This happens because the frame interval $[P_1, P_2]$ consumes too much of the total heat parameter space. Because of the large incremental heat value $\Delta q_i$, the subsequent heat values $q_i$ through $q_N$ tend to be pulled toward the key position $P_2$. In parallel with the heat values, the frame positions are also driven toward the key position. One consequence of the concentration in $\Delta q_i$ is that it propagates to the position $P_2$ when it tries to lower the next heat frame $P_1, P_2$ to the incremental heat value of $P_N$ should be assigned a value first cause the sum of the incremental heat value is around one. In practice, the portion of $P_N$ is not evaluated as part of the inherent frame or the interval $[P_1, P_2]$. Since it is the initial key frame of the interval $[P_1, P_2]$, it is evaluated with least parameter $q = 0$ at the interval and the position cannot be changed. Consequently, the interbetween position $P_2$ comes near to the position $P_N$, which belongs to the key frames and the motion speed gets slowed down relative to the previous motion.

Comparing intervals $[P_1, P_2]$ and $[P_2, P_3]$, the symmetry of the frame positions with respect to position $P_2$ is not maintained. For the same reason that $P_1$, the posture of the frame $\Delta 1$, is overly displaced from frame $\Delta 1$, frame $\Delta 1$ shows a large displacement from frame $\Delta 1$. Incremental heat layering is dependent upon the path it takes to produce the heat value. For instance, in the interval $[P_2, P_3]$, if the heat values are evaluated starting from the keyframe position $P_2$, up to frame $P_3$, frame $\Delta 1$ shows a large displacement from frame $\Delta 1$ and frame $\Delta 1$ displacement from frame $\Delta 1$ and the frames generated in both intervals will exhibit exact symmetry

Figure 6.14 illustrates the frame positions appearing after Lookahead Adjustment is applied to incremental Juza Spacing in Figure 6.14. Notice that the incremental time value of the first frame has been reduced as a result of Lookahead Adjustment. Since the position of the second frame with Incremental Error Spacing was overly displaced from the first frame, Lookahead Adjustment did look ahead to the second frame position and readjusted the incremental time values. The reduction led to a decrease of the Calculated Distance and thus the reduction of the Measured Distance. Note, however, this produced a slightly smaller distance between the first two frames than the expected mean distance d1. As a result of this readjustment, the sum of the incremental lens spacing falls below d₀₁, as evidenced by the long distance between frames #5 and #6. The underestimation can be corrected upon iterative application of Lookahead Adjustment as suggested in the section on Pseudacode Representation. Overall, the sum of the absolute distance error is reduced and the standard deviation from the mean value is greatly reduced, the inter-frame distance is closer to constant compared with Incremental Error Spacing. That is, the application of Lookahead Adjustment further aids in stabilizing the output speed, albeit a constant at this time.

Although Lookahead Adjustment further facilitates the control of the motion speed, there should be a way to let the user of Incremental Juza values to see in tenser proper choices between the last intervalrated frame in the first interval and the first frame in the last interval as described before. Figure 6.17 and Table 6.4 show the result of Averaging Adjustment applied to Lookahead Adjustment in Figure 6.15. The analog elongated inter-frame distances between $P_5$ and $P_6$ in Figure 6.16 has been redistributed through out the inter-frame distance in the interval $[P_1,P_6]$. The redistribution

correspends to picking all of the three positions in $[P_1, P_2]$ toward frame #5 at $T_p$. In parallel with the distance adjustment, the incremental knot values are redistributed while ensuring that they sum to one. By redistributing the distance moves over the entire interval, the transition from the last unconverted frame (i.e., frame #5 at $T_p$ in this case) to the first interconverted frame (i.e., frame #4 at $T_p$) becomes smoother than is the case with no

Averaging Adjustment: The smoothness near the key frames (frames #1, #4, #11 at positions $P_1$, $P_2$, $P_3$ in this case) is important since in our approach, the motion speeds in each pair of key frames are separately controlled and the frames are constrained to produce these positions operating multiple key frames.

Figure 6-19 illustrates an example in which the speed is controlled to vary between frames. The outer motion path as in Figure 6-17 is used to produce a total of 16 frames in Figure 6-19. Three discrete speed values, 8, 1, 8 are assigned for the key frames at $P_1$, $P_{int}$ and $P_3$, respectively. One can specify 4, 0, 4 for the more speed value however, since only the relative values of the speed matter. The inner three clusters of the frames at $P_2$ is organized to be right since shorter than the two near the frames at $P_3$ and $P_{int}$. To produce a gradual speed change, the inbetween frames are set to take the speed values calculated by linear interpolation, such that the speed of the $O\!-\!th$ inbetween frame in the interval $[P_1, P_2]$ is given by

$$s_i = \text{speed of the moving key frame} \times \frac{(O-i)}{\text{number of frames}}$$

$$+ \text{speed of the ending key frame} \times \left(1 - \frac{(O-i)}{\text{number of frames}}\right)$$

$$= 1 \cdot \frac{O-i}{O} + 8 \cdot \left(1 - \frac{O-i}{O}\right)$$

where

the frame number i in an integer ranging from 1 to 6.

It is important to note that the temporal aspect is completely separated from the spatial aspect in our approach. The speed profile can be independently generated from the motion path. For instance, the motion speed is controlled by a linear interpolation method while the motion path is produced by the cardinal spline Incremental Knot Spacing followed by Lockianhead Adjustment and Averaging. Adjustments were used to generate Figure 4-18, in which the motion draws an explicit slowdown near the key frame positioned at $P_4$ while it speeds up near the key frames positioned at $P_2$ and $P_3$. A measure of error, the difference between the Calculated Distance and the Measured Distance, sums to 56 in cm for rows in Table 4-5. Most of the errors are distance errors near frame #1 positioned at $P_1$ and this can be verified by comparing the speed profile curve and the measured speed in Figure 4-19. In the figure, the speed profile curve locations (those) with respect to the frame sequence number because of the linear speed interpolation scheme given above. However, the measured speed shows some deviation from the speed profile curve near the first frame and the error comes from the fact that the tangential magnitude is underestimated near the first frame, causing large incremental joint values.

Figure 4-20 shows the relationship between the frame sequence number and the knot values in Figure 4-18. As was mentioned previously, in our control scheme, the motion speed is controlled by adjusting the knot parameter value t corresponding to the temporal knot number. Roughly speaking, the inter-frame distance is proportional to the magnitude of the incremental knot spacing. In view of the fact that the slope at a certain frame in Figure 4-20 can represent the incremental knot spacing, the slope near frame #3 represents the incremental knot spacing, the slope near frame #9 is small compared with that of frame #1. Therefore, the deceleration

of speed in Figure 4-24 is driven by the decrease in the incremental heat spacing in Figure 4-22.

Figures 4-25 and 4-22 illustrate the position of inbetween frames and the measured speed values when the number of inbetween frames is double that of Figure 4-19. The same median path and method for the motion speed control was used in both figures. As can be seen, the inter-frame distances ensure including the frames near the first frame tend to be smaller and thus the measured speed tends to approach the speed profile curve. Another measure of error, the sum of absolute distance errors, is much reduced compared with Table 4-8. In fact, the error approaches zero as the number of inbetween frames increases.

Figures 4-26 through 4-27 illustrate how rapidly the errors can be reduced in one motion control method by increasing frame numbers with different median paths shown in Figure 4-23. The results in the figures are based on a speed profile with a constant speed, and therefore constant inter-frame distances. To compare the behavior of our speed control scheme in the longer interval ($F_a$-$F_a$) and the shorter interval ($F_b$-$F_c$), Figure 4-23 is made to be unsymmetrical. As a measure of error, the sum of the incremental knot values is used in Figures 4-26 and 4-27. If Averaging Achievement is employed, the sum is identically zero. Therefore, the sum can be used to estimate how accurately Incremental Knot Spacing and Lookahead Adjustment can approximate the given constant speed profile. As long as it approaches one, the individual knot spacing and corresponding inter-frame distances can be said to be measured properly. As another measure of error, the sum of the Distance Error is used in figures 4-26 through 4-27. Part of the reason why it is used instead of the sum of the absolute Distance Error is that it can be a reference of whether the knot Calculated Distance exceeds the Measured

Distance by including the sign. This error should approach zero as the number of abscissas frozen increases, so that the Calculated Distance can properly approximate the Measured Distance.

Depending upon the interpolation scheme of the motion path, the performance of our speed control method causes it a relatively small number of in-between frames are incorporated. The transition occurs since each interpolation scheme has a different mathematical representation of the motion path and the difference leads to variation of required magnitude replicated in our speed control method. Figures 4-24 and 4-25 illustrate the sum of incremental knot values in the intervals $[P_0, P_9]$ and $[P_9, P_{16}]$, respectively. The values are taken when Incremental Knot Spacing followed by Localized Adjustment are applied to three different interpolating schemes for the motion path. At a small number of in-between frames, the sum of the incremental knot values of the Type II free form spline exceeds one, while that of the cardinal spline and the B-spline falls below one. The cardinal spline falls furthest away from one than the B-spline. This means that the Type II free form spline and the cardinal spline have a greater curvature in this region, causing the tangential magnitude to vary rapidly enough for a fixed request to be in properly represent the entire range of tangential changes happening between a pair of in-between frames. Note that the tangential magnitude cannot be determined directly from the appearance of the curve shape in Figure 4-23. It can be evaluated only by mathematical representation, as mentioned earlier in the interval $[P_0, P_9]$, however, only the sum of the incremental knot values of the Type II free form spline exhibits convergence to one before the number of in-between frames reaches five, so can be seen in Figure 4-25. In view of the convergence, this equal convergence compliant with the general behavior of the G² free form

splicer explained in Chapter 3. G7 has been spline has relatively constant curvature in the turnover interval so that a fixed tangent can be effectively used to estimate the tangential magnitude in each substantial scale up of the motion path between a pair of key frames. Overall, the sum of the incremental keys values show rapid convergence to one unit relatively low inbetween frames.

The distance series in figures 4-56 and 4-57 are the relative magnitudes of the distance when the distance between the keyframe pointers $P_3$ and $P_4$ is set to one. Incremental Knot Spacing, Incremental Adjustment and Averaging Adjustment were used to generate the inbetween Points leading to the figures. As the number of inbetween frames increases, the distance series approach zero so that the calculated distance makes little deference to the measured distance. In figure 4-56, the sum of distance errors of Type U line knot splits is to the negative side, while others are positive with a low number of inbetween frames. This means that the calculated distance is smaller than the measured distance.

It is interesting to compare the behavior of the curves with respect to the vary of the incremental knot values with that of the curves with respect to the sum of distance errors. For instance, if the sum of the incremental knot spacing of Type R has lower spline at 0.61 in greater than one, it is less than zero in Figure 4-59. Other curve tolerate exhibit similar behavior. The same is true for the Figures 4-55 and 4-57. As the sum of the incremental knot spacing becomes greater than one, the sum of distance errors tends to be less than zero. The principal reason for such behavior is that Averaging Adjustment counteracts the error in the sum of the incremental knot values. If the sum of the incremental errors is large, Averaging Adjustment makes the sum into one by reducing the calculated distance to the point where the

sum of the colocated distance is less than for the rest of the measured distance. This is true in the case of a small number of the otherwise frozen itlerants, at the number of frozen increases, the counteracting property of Averaging Adjustment is small enough to be ignored and the distance errors approach zero. In conclusion, Incremental Knot Spacing combined with Lookahead Adjustment and Averaging Adjustment proves fairly effective in terms of the associated errors in the number of itlerants frozen increases.

## Summary

This chapter has been devoted to the control of motion speed in transition. In connection with our anomalies system, we developed a method to control the motion speed of a single state so that the displayed frame sequences do not exhibit unnecessary acceleration, or even random motion dynamics. Close tethered to the problem relating our to visualised with one of the motion is controlled such that a smooth transition between frames is ensured. In this chapter, we have presented a flexible speed control scheme for that purpose.

The basic building block of our method is Incremental Knot Spacing, which dynamically adjusts the parameter values of adjacent inbetween frames, depending upon various parameters: Lookahead Adjustment is employed to avoid relatively large errors in the incremental knot values. Averaging Adjustment was introduced to secure a controlled transition from the last inbetweened frame to the next key frame.

Figure 6-12. The positions of inbetween frames in a truncated cardinal spline with a tension value of one. Incremental lever spacing $\Delta s$ is set to a uniform value of 0.1. The positional vectors $P_{s1}P_{s2}$ and $P_{s3}$ represents the key frame positions. Within regions $[P_{s1}P_{s2}]$, through $P_{s3}$ are the positions of the inbetween frames produced by evaluating the cardinal spline at corresponding knot values ($t = 0$ through $t = 1$).

Figure 4-13: The positions of sixteeen knotats on a line spline with midpoint knot parameter values of $u = 1$, and $v = 0$ (e.g. .025) (incremented knot spacing as is set to a constant value of 0.1. The positional vectors $P_1$, $P_{11}$ and $P_{16}$ represent the key knots. Within region $[P_1, P_{11}]$, $P_2$, through $P_9$ are the positions at the sixteeen knotats produced by evaluating free form spline at corresponding knot values ($u = 0$ through $v = 9$).

Figure 6-14 Interhemural frame positions with constant laser spacing.

| Frame Number | Longitudinal Magnitude | Isr Knot Value | Measured Distance |
|---|---|---|---|
| 1 | 160 | .200 | 41 |
| 2 | 264 | .200 | 79 |
| 3 | 394 | .200 | 81 |
| 4 | 506 | .200 | 71 |
| 5 | 384 | .200 | 41 |

Table 6-1 Inter-frame distances of Figure 6-16

Sum of incremental knot value = 1.000

Standard deviation of measured distance = 16.8

Figure 6-10: Inferrence frame positions upon application of Incremental Knot Spacing

| Frame Number | Tangential Magnitude | Int Knot Value | Calculated Distance | Measured Distance | Distance Error |
|---|---|---|---|---|---|
| 1 | 168 | .361 | 41 | 322 | 41 |
| 2 | 388 | 157 | 41 | 63 | |
| 3 | 603 | 150 | 41 | 58 | 3 |
| 4 | 809 | 179 | 41 | 53 | 0 |
| 5 | 142 | 201 | 41 | 20 | 5 |

Table 6-2. Error analysis of Figure 6-10

Sum of concentrical knot value = 1.110, Sum of absolute distance error = 90, Standard deviation of measured distance = 27.4.

Figure 6-9c Intertwist linear position upon application of Lockabout Adjustment

| Frame Number | Tangential Magnitude | Arc Axis Value | Calculated Distance | Measured Distance | Distance Error |
|---|---|---|---|---|---|
| 1 | 140 | .200 | 37 | 49 | -12 |
| 2 | 309 | .169 | 52 | 48 | 4 |
| 3 | 309 | .131 | 39 | 41 | -2 |
| 4 | 605 | .059 | 44 | 41 | 3 |
| 5 | 390 | .250 | 73 | 73 | 1 |

Table 6-3 Error analysis of Figure 6-9c
Sum of incremental load value = 317,
Sum of absolute distance error = 18,
Standard deviation of measured distance = 7.5

Figure 4-17 Inbetween frame positions upon application of Averaging Adjustment

| Filter Number | Temperature Magnitude | Bar Rate Value | Calculated Distance | Measured Distance | Distance Variance |
|---|---|---|---|---|---|
| 1 | 160 | 349 | 60 | 96 | -10 |
| 2 | | 280 | | 66 | |
| 3 | | | | 67 | |
| 4 | 395 | 220 | 66 | 67 | -1 |
| 5 | 310 | 290 | 71 | 52 | 21 |

Table 4-4. Error analysis of Figure 4-17

Sum of incremental bcar value = 1 045.
Sum of absolute distance error = 55.
Standard deviation of measured distance = 7.6

Figure 4-18: Linear speed interpolation with 8 interznent frames in intervale $[P_4P_5]$ and $[P_5P_6]$.

| Frame Number | Speed | Tangential Magnitude | Iss Rate Value | Calculated Distance | Measured Distance | Frame Error |
|---|---|---|---|---|---|---|
| 1 | 8.0 | 163 | 267 | 40 | 55 | -18 |
| 2 | 767 | 312 | 164 | 53 | 60 | -7 |
| 3 | 6.99 | 397 | 139 | 81 | 85 | -4 |
| 4 | 5.34 | 605 | 111 | 85 | 44 | 1 |
| 5 | 4.36 | 734 | 106 | 97 | 93 | 2 |
| 6 | 3.61 | 916 | 108 | 107 | 113 | 3 |
| 7 | 3.78 | 289 | 064 | 29 | 33 | 3 |
| 8 | 3.78 | 289 | 064 | 29 | 33 | 3 |

Table 4-5: Error analysis of Figure 4-18: Sum of absolute distance error = 34

Figure 4-19 Comparison of the speed profile curve with measured speed which is the speed profile curve reconstructed from the perceived tolerance frames in Figure 4-20.



Figure 4-20 Relation between the frame sequence number and the normalized knot parameter value used in controlling the motion dynamics features in Figure 4-19.

Figure 4.21: Linear speed interpolation with 16 interbreath frames in intervals $[F_1,F_2]$ and $[F_2,F_{39}]$.



Figure 4.22: Comparison of the speed profile curve with measured speed which is the speed profile curve reconstructed from the generated inbetween frames in Figure 4.21.

| Frame Number | Distant Speed | Integrated Magnitude | Sin Root Value | Calculated Distance | Invariant Distance | Distance Score |
|---|---|---|---|---|---|---|
| 8 | 4.92 | 1.62 | .179 | 25 | 31 | -6 |
| 9 | 7.24 | 2.17 | .238 | 29 | 31 | -4 |
| 7 | 3.22 | 1.02 | .065 | 27 | 31 | -2 |
| 6 | 6.49 | 3.75 | .071 | 37 | 38 | -1 |
| 5 | 4.25 | .999 | .064 | 28 | 28 | 0 |
| 4 | 5.51 | .657 | .126 | 34 | 33 | 1 |
| 3 | 7.58 | .466 | .055 | 22 | 22 | 0 |
| 2 | 4.94 | .265 | .022 | 31 | 30 | 1 |
| 9 | 4.50 | .759 | .190 | 17 | 16 | 1 |
| 11 | 5.43 | .442 | .066 | 20 | 18 | 2 |
| 10 | 3.58 | .567 | .047 | 39 | 37 | 2 |
| 12 | 3.18 | .396 | .349 | 14 | 11 | 3 |
| 13 | 3.75 | .261 | .065 | 12 | 9 | 3 |
| 14 | 7.08 | .542 | .075 | 18 | 15 | 3 |
| 15 | 1.90 | .261 | .002 | 9 | 7 | 1 |

Table 6.4  Error analysis of Figure 6.21
Sum of absolute distance error = 30
Sum of individual lengths = 306

Figure 4.23: Different interpolation methods for the motion path generation. Position of key frames are represented by $P_i$. Figures 4-21 through 4-27 are based on the subsequent frames generated by these interpolation methods. The gray dots represent the position of the key frames such as end points for the interpolatory curves in the intervals $[P_0P_1]$ and $[P_3P_4]$.

Figure 6.24. Sum of discounted level values in the interval $[P_1, P_2]$



Figure 6.25. Sum of discounted level values in the interval $[P_1, P_2]$

202



Figure 4-26 Sum of the distance errors in the interval $[P_2,P_3]$



Figure 4-27 Sum of the distance errors in the interval $[P_3,P_4]$

# CHAPTER 5
## APPLICATION: SCRIPTS FOR PROTEIN FOLDING

### Introduction

This chapter describes an animation system as an application of the preceding chapters. More specifically, a prototype system is designed for the simulation of the protein folding. The free form splines developed in Chapter 3 form the basic building blocks of the animation system. First, it is used for producing a curve representing the backbone conformation, and second, it is used for the generation of a motion path of individual backbone atoms. The speed control techniques developed in Chapter 4 are used to generate a smooth motion transition between frames. In addition, inbetweening techniques discussed at Chapter 4 are introduced, for the first time in the field of molecular graphics, at the context of the backbone animation.

To facilitate exploration, a general background of protein folding in biochemistry is presented. Terms appearing in subsequent sections will be defined briefly in this section. Following a summary of previous work in the area of molecular graphics, the problems associated with current molecular animation systems will be addressed. Based on observation of these problems, fundamental concepts which differentiate our prototype system from other conventional approaches are explained. Subsequently, design topics involved in the implementation of the prototype system are depicted. Finally, some of the frames produced from our prototype system are illustrated as experimental results.

196

## Protein Geometry

Biochemistry is the study of the molecular basis of life. For instance, the discovery of the double-helical structure of deoxyribonucleic acid (DNA), the elucidation of the flow of information from gene to protein, the unraveling of the energy-conversion mechanisms, the determination of the three-dimensional structure and mechanisms of action of many protein molecules, are some of the outstanding achievements of biochemistry.

In virtually all biological processes, proteins play crucial roles in enzymatic catalysis, transport and storage of ions and small molecules, coordinated motion as a muscle, mechanical support of skin and bone, immune protection, generation and transmission of nerve impulse and control of growth and differentiation.

Amino acids are the basic structural units of proteins. An α-amino acid consists of an amino group, a carboxyl group, a hydrogen atom and a derivative R group bonded to an α-carbon atom. An R group is called a side chain and distinct amino acids result from the variation of its component atoms in proteins, the α-carboxyl group of one amino acid is joined to the α-amino group of another amino acid by a *peptide bond*. Figure 9.1 illustrates the formation of a peptide bond. Residues 1 and 2 can be regarded to form a polypeptide chain. An amino acid unit in a polypeptide is called a *residue*. By convention, the amino group of the N. terminus, is taken to be the beginning of a polypeptide chain. A polypeptide chain consists of a regularly repeating part, called the *main chain* (regardless of nitrogen, alpha carbon, carbonyl carbon), and a variable part, comprising the side chain or chains. The main chain is sometimes termed the *backbone*.

A remarkable characteristic of proteins is that they have well-defined three-dimensional structures. the rigid and planar peptide unit. Figure 5.6 shows a pair of peptide planes. The hydrogen of the substituted amino group is almost always opposite the oxygen of the carbonyl group. No freedom of rotation about the bond between the (carbonyl) carbon and the nitrogen atom of the peptide unit exists, since this peptide linkage has a partial double-bond character. In contrast, the link between the alpha carbon atom and the carbonyl carbon atom is a pure single bond. and the bond between the alpha carbon atom and the peptide nitrogen atom is also a pure single bond. Therefore, there is a large degree of rotational freedom about these bonds on either side of the polypeptide unit: There is, relative to the peptide bond, free rotation about the bond between the alpha carbon and nitrogen atom. The same is true for the alpha carbon to carbonyl carbon bond. The rotational angles about these bonds are designated phi and psi, respectively. The conformation of the main chain of the polypeptide is determined when the phi and psi angles for each amino acid are defined.

As a result of this freedom, some three-dimensional patterns or protein structure can be found. Alternately these patterns are called secondary structures as a conceptualization of a sequence of residues. The alpha helix is a stablize structure. As shown in Figure 5.9, the inner part of the rod is made up of tightly coiled polypeptide backbones, while the side chains extend outward. What stabilizes the structure is the hydrogen bond. The alpha helix is stabilized by hydrogen bond between the NH and the CO groups of the backbone. The CO group of each amino acid is hydrogen bonded to the NH group of the amino acid that is located four residues ahead in the linear sequence.

The β (phenol) sheet differs markedly from the α helix in that it is a sheet rather than a rod. Figure 5.4 illustrates a β sheet structure. A polypeptide chain in the β sheet is fully extended, rather than being tightly coiled as in the α helix. Some other secondary structures (GlnH, ChH) such as β-turns, which are required for the directional change of the β strand chain, known as hairpin bends or α helix to β sheet, also form a short directional change from the α helix to other structures, can also be found in nature. C. N. Ramachandran and V. Sasisekharan at the University of Madras, India, showed that anthocarilic interactions between late-carbon atoms or amino acid side chains with some of the polypeptide backbone make only three repeated conformations favorable; the right-handed alpha helix, beta strands and the left-handed triple collagen helix. Alternatives produce turns, loops and random coils.

Protein engineers also define the tertiary structure (3oHL) of molecules using the secondary structure as a basic building block. They expand the secondary structure into a pattern called the domain (R₂H1), where alpha helices and beta sheets form surface shapes such as barrels by parallel and anti-parallel combinations. In this respect, the recognition of hairloop shape is essential for the recognition of a higher level structure.

In an alpha helix, each residue is rotated 100 degrees relative to the previous one, and it is translated along the axis by about 1.5 angstroms (each turn containing 3.6 residues). Every carbonyl oxygen points upward and receives a hydrogen bond from a downward-pointing amino nitrogen of the fourth residue away from it. In this polar structure, the polypeptide backbone is intimately extended and each residue is rotated 100 degrees with respect to the previous residue. The strands pack side by side to form a β sheet, in which each amide nitrogen donates a hydrogen bond to a carbonyl oxygen in the

adjacent annual and such carboxyl oxygen moieties such a bond. Acceptable pts
and pss bond angles which correspond to those in sheet conformations can be
determined by extending a linear/random plot.

Polypeptide chains interact with their environments to fold up into
discrete, highly organised and tightly packed three-dimensional structures.
Diverse coiled shapes confer on proteins their powers of recognition and
selectivity, catalysis, and specificity as structure-forming elements.
Determining the conformations of proteins is therefore critical to
understanding their roles in biological systems and biochemical pathways.
Most of the three-dimensional structural data obtained for biological
macromolecules comes from X-ray diffraction experiments or the electron
microscope. X-ray diffraction analysis requires that molecules be ordered into
fibres, sheets or crystalline arrays, however, not all macromolecules crystallise
to form such ordered structures. Brittkess microscopes resolve to about five
angstroms, which is less than the accuracy required to determine three
structures. A simple method for predicting how a protein will fold is clearly
needed.

The initial breakthrough in the study of protein folding occurred in
1957 when Christian B Anfinsen unfolded and refolded purified ribonuclease
in the absence of cellular components. This study established that all of the
information needed to determine the final conformation of a protein can
reside in the polypeptide itself. Here the sequence of residues determines the
folding of proteins, and hence their conformations, has remained a major
unsolved question in biochemistry.

The term conformation or general term encompasses the spatial
arrangement of a molecule as determined by rotations about the single bonds.
In the specific case of the polypeptide chain, conformation describes the

overall spatial organization. The description of polypeptide conformation involves the specification of bond length, bond angles and angles of internal rotation about the single bonds.

At the beginning of the list of the chain, residues in the sequence interact with each other and the solvent to cause the chain to fold up. These interactions are almost solely dependent on the nature of the side chains which characterizes the twenty amino acids. Charged residues, such as aspartic acid, glutamic acid, and tyrosine, are located preferentially at the protein surface, where they can interact with water; residues in the interior are clearly packed, with the solvent relatively excluded. Burying of the hydrophobic groups and surfaces is a major source of the protein stabilization. For hydrophobic amino acids, interactions exclude both van paar and hydrogen bonds. Hydrophobic amino acids are probably critical in the forming of solvent-excluding protein interiors. Within proteins, the conformation of the amino acid side chains are not retained, but represents low energy state.

Recent progress on polypeptide chain folding has come from work in three areas:

1. Investigations of the determination of folding interactions using hydrogen exchange and nuclear magnetic resonance (NMR) techniques.

2. Pinpointing the residues in polypeptides that carry the information that determines the folding pathway by the use of polypeptide chains with amino acid substitutions at known positions.

3. Analysis of traps within cells between the nascent polypeptide chain polymerized sequentially on the ribosome and the appearance of the native protein. For reported proteins, there is a complex state of

interactions of the nascent chain with cellular components, including helper proteins called molecular chaperones.

All these approaches are used to further elucidate the possible folding pathways.

Aside from the characterization of the role of the amino acids or protein folding, there have been computer-based approaches in determining the folding intermediates, and thus the folding pathways. Energy minimization is a predominant simulation technique to determine a reaction pathway, since large movements which involve advanced displacements of many atoms or those requiring the crossing of energy barriers cannot be investigated with molecular dynamics [ile77, fis91].

In [l7aa91, lar91], energy minimization is used extensively as a simulation tool for the inspection of the conformational change occurring at large molecules. We note are derived from cholesterol and their formation is the principal cause of cholesterol degradation and elimination. In [lar91], chalic acid (a kind of bile acid) layers are built on a screen by stacking individual chalic and molecules together. Subsequent energy minimization did clarify the change at the tilde and shape and orientation of channels peculiar to the micelle. In [l7aa91], adipisitic molecules are joined together and these macromolecular change investigated. Depending upon the peak positions, various patterned structures were derived so that one solid arouse that a regulatory function was ruling the conformation. Since energy minimization destroys the protein conformation in the direction of a more stabilized form, it is amenable to assure that the folding process can best be simulated with it.

Therefore, the subject conformations from energy minimization calculations are not directly an input feature in the simulation. However, the

mass display of sequential frames from energy trajectories does not result in smooth conformational change. The degree of conformational change from one frame to another varies widely. For instance, in one frame, some portion of a molecule may deviate from the previous frame by a wide degree of displacement while in next frame, the same portion remains relatively fixed. Energy minimization generates conformational data relying solely on energy levels, and spatial distance in the conformational data cannot be reflected in its calculation.

## Laboratory Reviews

### Molecular Graphics

Since Cyrus Levinthal's first attempt to draw molecules on a computer in 1966 [Levin], molecular graphics or a field of computer graphics has undergone extensive research and improvement. Traditionally, static and tedious mechanical reconstruction of molecules in limited laboratory space has been the impetus for the proliferation for computer representation. In addition, computer representation permits rotation and translation of a three-dimensional representation to aid in the recognition and understanding of a three-dimensional molecular structure [Schol9].

Molecular graphics is an ever-expanding field with immediate practical impact on the design and production of structures for useful, pharmacologically-active proteins. Recent sophisticated demands require the control of specific protein conformations such as bond lengths, bond angles, molecular surface areas, molecular volumes, or hydrogen operation on atom and volumes with a simple interface device such as a mouse [Stel6, Zupfe].

Individual atoms and the bonding relations between atoms fall into a few general classes so the representation techniques which have been applied to date. Among these computer graphical representations, the skeletal model, ball-and-stick model and space-filling model have their mechanical counterparts engendered by truncated plastic materials.

The skeletal model (or endpoint model [&309] is based on Kendrew-Watson's physical model of a molecule, typically with a scale of 1 Angstrom unit being equal to 20 mm and with a wire diameter of 5mm. This model shows only the molecular framework, as a collection of lines joining the atoms' positions of bonds and their termini imply atomic position. Advantages of this model are elegantly and the capability of rapid imaging of large molecular (e.g., these reductions where the number of atoms exceeds 1,000). The skeletal model permits the viewer to determine bond lengths, angles and dihedral angles [&309, &e91, &a98] with ease.

Devised by Corey, Pauling, and Koltun, the CPK model (space-filling model or spherical model) illustrates molecular surface shapes; atom radius shapes determine each atomic's background function as internal-inclusion interactions [&309]. While an atom is a spectrum excluded volume where volume is specified by a probability distribution, experiments show that spheres can approximate atomic volume and shape [Che50]. The van der Waals contact distance [Me97] determines the atomic radius in the CPK model. This distance, defined as the length at which a repulsive force begins to exert an influence when two atoms are drawn as such atomic, a larger or smaller, unit touches. Two non-bonded molecules such as bonding at vertically complementary atoms, van der Waals interaction at caused by temporal asymmetry in electron charge attracts atoms together. Therefore, the CPK model treats the van der Waals radius as the size of an individual atom and represents the

molecules at the centre of spheres (Motel8, Net5??). Nevertheless, the bonds between water are invisible in the spherical model. Because the bonds are the overlapping region of the electron density that becomes invisible in the spherical model, they are bushed under spherical surfaces and the relative geometry between the atoms is hard to distinguish.

The dot sphere model uses dots to depict atomic surfaces or outlines of a constant electron density. The ball-and-stick model induces the size of an atom relative to bond length to make the inter-atomic bonds visible. Varying the radius of the connection link in the ball and stick model imparts a sense of relative depth between two atoms as in a perspective view transformation. Stereographic views pairs of this model (SylPH4) allows depth perception of interesting lines.

The interaction between the protein surface and a water molecule offers another modeling issue. The colour scientist surface model (Connll, Lee61, Maxel6) as a variation of the spherical model, makes it easy to visualise the entomology of hydrophobic and hydrophilic reactions. This model displays Van der Waals surface areas of an atom that is accessible to the spherical surface of a water molecule.

With the growing popularity of the CPK model, various efforts at computer graphics area have been aimed at reducing their enactment involved in its use.

Polygonal approximations of a sphere is a standard surface normal vector to speed up hidden surface removal (Seg68, Net5??). To avoid black edges on polygonal facets, the number of facets must be increased enormously based on the scaling of color, generalisation of the shading function for a sphere SixM4 was applied on a scan line basis to reduce the black effects by elimination of polygonal approximations. This approxi-

shading algorithm is implemented with a dedicated display buffer [Sei84] to speed up the rendering of spheres at large resolution. In an attempt to reduce the complexity and thereby to reduce the time constraints required to draw intersecting spheres, a model of molecular structures with non-intersecting spheres has also been investigated [Sta69].

Rather than subdividing a sphere into a collection of polygons, Porter [Por78] treated a sphere as a primitive and explored the fact that a sphere projects into a circle regardless of the position of viewpoint. Porter applied Bresenham's image space circle generation algorithm [Bre77] to avoid square-root calculations. However, this method suffers from aliasing [Cro77, Cro81] by drawing jagged edges. Such aliasing effects can be greatly reduced by the use of post-aliasing lookup tables [Wan83]. If one simplifies a sphere into a circle [Bra75, Max75, Max83, Knu83], a conventional hidden line removal method can be used to calculate visible area. In this method, a circle of filled circle will approximate a sphere, but the model lacks visual reality. In the case of non-overlapping spheres, an object space pruning algorithm [Bra81] is possible. By enclosing the surface which are declared hidden or prevent removal stages, the algorithm maximizes the complexity linearly proportional to the number of spherical objects.

Equally as importance to the spherical model is the ribbon model, which has been used to visualize a higher level of conceptualization, especially in protein engineering. Because the amount of graphical information in a globular protein is very large, researchers tend to lose interest if they are confronted with the full details of the individual atoms represented by a spherical model. Rather, they need to interpret local distribution of atoms in terms of global secondary structure, which are regular arrangements of the sequence of residues. For instance, biochemists are interested in whether they

can find a particular amino acid ([Ja69] or a particular location or a long residue sequence. In this respect, the ribbon model, as a framework for the recognition of higher level structures, such as the tertiary structures, plays an important role.

Although protein engineers may fit together several side chains during the design process, their ultimate interest lies in the backbone conformations, since the stability of the molecule is determined by the spatial pattern of such conformations. The ribbon model [Ric81, Ric85] in this context portrays a molecule at multiple threads approximating the lines joining backbone atoms. Although the backbone atoms are defined to include all the nitrogen, alpha carbon and carbonyl carbon atoms, the line joining the adjacent alpha carbons has been treated as the backbone line alternatively. Figure 8-5 illustrates the alternative backbone line. Figure 8-6 illustrates a ribbon model representation of ribbon conformational structures. Instead of a single thread, multiple threads are drawn in parallel to help the viewer perceive three-dimensional depth on a two-dimensional screen. Various patterns of naturally stable secondary structures such as heles, sheet, and turn nets can be visualized with ease using the ribbon model.

Overall, the exploration of molecular graphics to date has mainly concentrated on a static display of a molecule and the use of models varied to suit the diverse needs of the protein engineers. Depending on what aspect of a molecule should be emphasized and maximized, various models have been proposed. Nevertheless, the modeling aspects with regard to the animation of molecules have not yet been refined. In the following section, we will examine and analyze how previous molecular animation systems have been developed.

## Molecular Animation

With the growing interest in molecular graphics, the obvious extension of existing means of computer representation is an animated graphic display of the molecular motion. Feldmann and Levitt (Ref68), at a 1980 film of the molecular dynamics of bovine pancreatic trypsin inhibitor (BPTI), demonstrated the power of computer animation as a tool for the researcher. Fundamental concepts were easy to grasp using visual effects involving motion. The technique of computer animation shows promise in applications from non-erasure tubular surgery to the prediction of structure for enzymes, and to the scattered display of such molecular processes as DNA unraveling or protein folding/unfolding synthesis.

In the decade since these earliest efforts at computer animation of the motions of a molecule, the computer-aided design, manipulation and display of molecular structures have grown into a new hundred million dollar-a-year business. However, most of the existing molecular graphics packages put great emphasis on enhancing the photorealism of a static view of molecules, while ignoring the potential for animated display of molecular motion.

These graphics systems which do attempt the animation of molecular dynamics data to simulate molecular motion share a number of common characteristics as follows:

Foremost of these is that they require the use of a mainframe computer, since the generation of realistic images can be prohibitively expensive. Because the computer representation of a spherical or space-filling model requires an enormous amount of CPU time, the simpler wire-frame model has proven more useful in all attempts to date at molecular

animation. The wireframe model or an animation tool (Delft, Zajfe, Tschöl) permits rapid imaging of large molecules, but has the disadvantage of each distracting visual effects as flipping side chains, abrupt changes in angles and disjointed motion. As a result, viewers are easily distracted from the global changes in conformation. Second, conventional concepts of animation as a series of sequential displays of a frame have proven of limited use when applied in molecular dynamics. Third, most of the approaches to date have concentrated on interesting interactivity (Delft, Zajfe, Tschöl), so that a user could control animation speed and viewpoint, and freeze the animation screen in order to analyse individual conformations. Such controllability and the added capacity to analyse a static scene are meaningful in that they help recognize the three-dimensional structure of the molecules, but are useful only if a proper model for animation is employed.

To summarise the problem, currently available animation systems such as the Winchester Graphics System (Tschöl) or RIUGIG, (Delft) have shown too much emphasis on the implementation of the controllabilities. The molecular model predominant in the previous animation system was the wireframe model, which continually produced such visual anomalies as high-frequency motion of flipping side chains and abrupt motion discontinuities. Essentially, the wire frame model is not suitable for animation. In addition, the conventional systems could not show smooth movements, since a series of halting intermediaries, or so, is displayed without any consideration of the possibility of generating in-between frames.

218



Figure 3-7. Peptide bond between residues. The atom types are, H for hydrogen, C for carbon, O for oxygen, N for nitrogen. $C_\alpha$ for alpha carbon. The R group is a collection of atoms which characterize the given residue.

Figure 5-2. Definition of α and φ refers to rotances about the $C_β$-C single bond; φ refers to rotances about the $C_α$-N single bond. In a fully stretched-out polypeptide chain, φ = ψ = 180 degrees

hydrogen bond

1.0 Å rise
100 degrees rotation

9 Å

Figure 9-9. A schematic diagram of a right-handed α-helix. Hydrogen bonds between NH and CO groups stabilize the helix. In the α-helix, the CO group of residue n is the hydrogen bonded to the NH group of residue (n + 4).

Figure 9.4 A schematic diagram of the anti-parallel β pleated sheet. Adjacent strands run in opposite directions. Hydrogen bonds between NH and CO groups of adjacent strands stabilise the structure. The side chains R(i) are above and below the plane of the sheet.

In the figure: labels include "hydrogen bond" and "direction of the strand".

an alternative
backbone curve
the line segments
joining alpha carbons

a backbone curve
reconstruction of
the line segments
joining the backbone atoms

Figure S/i: A backbone conformation. Atom types are, N: nitrogen, C: carbon, O: oxygen, H: hydrogen. $C_\alpha$: alpha carbon. The T group represents a collection of atoms which characterize a given residue.

Figure 5-6 A schematic diagram of the protein secondary structures represented in the ribbon model. A coil as relatively tightly wound compared with a turn.

## System Overview

Two specific aspects of our statement system should greatly improve the quality of the correlational molecular statement systems. These are,

a) Modified use of the ribbon model

b) Introduction of key-frame inferencing

The $G^2$ free forms developed in Chapter 3 makes up the theoretical foundation for the modified ribbon model which will be discussed subsequently. The method to establish the ribbon model developed in Chapter 4 will be used as key-frame inferencing at our statement system.

### The Solid Backbone Model

Although the ribbon model has been used to generate a main view of protein residues, the shortcomings of the model are used as an extraction tool, can be summarized as follows:

1. When extracted, the multiple parallel threads making up the ribbon results in a hopping motion which distracts viewer's attention from the conformational change of the backbone.

2. The original purpose of drawing multiple threads at the ribbon model was to enhance depth perception as a main view. But the threads enhance depth perception only when a viewer has enough time to examine the width and twist of the ribbon. In a dynamic display of frames, however, the time between the frames is very limited since the frames should be displayed at a rate of about thirty frames per second in order to produce the visual illusion of animation.

3 In the ribbon model, the curve representing a sequence of the backbone atoms does not necessarily pass through the backbone atoms themselves.

In order to resolve these problems involved in the animation of ribbon model, we define the *solid backbone model* as a proper tool for the interaction as follows:

1 The solid backbone model is a combination of cylinders with their axes joining the backbone atoms. Thus the backbone curve appears similar to a tube or a cable.

2 The solid backbone model passes through the backbone atoms exactly.

Because it is defined as a three-dimensional solid (i.e., combination of the cylinders), the clipping surface of the multiple threads during the animation can be removed. Meanwhile, the three-dimensional depth of the protein conformation can be perceived through the shading on the surface of the cylinder. The shading is a more natural way to recognize a three-dimensional structure and it allows immediate recognition of the backbone conformation during animation. In a sense, the model is also suitable for the static view of the backbone conformations, since the only reason the ribbon model is used is the representation of the protein backbone in the depth perception.

Various illumination techniques for shading can be found in [Fol82, Rog85]. Usually, the illumination models are applied in a polygonal surface using a constant normal of the polygon. In the solid backbone model, a cylinder is made up of patches of small polygons and the number of polygons approximating the cylinder can be chosen arbitrarily. Although the image of a cylinder becomes close to a real cylinder as the number of the polygons increases, the computations involved in the calculation of the intensity of a

pixel becomes costly also. In animation, however, it is generally believed that a rough shading will suffice to represent an object, since viewers tend to ignore the relative details of the shading while the image changes.

Another technique involved in improving the three-dimensional perception of a solid backbone model is the use of hidden surface removal algorithms. These algorithms attempt to determine the faces or surfaces that are visible or invisible to an observer located at a specific point in space. The fundamental assumption made in these algorithms is that, the farther an object is from the viewpoint, the more likely the object is to be totally or partially obscured by one closer to the viewpoint.

Various solutions, such as a zone line z-buffer algorithm (ABA6, SUTE), visible surface ray tracing (WOAN, CLAV), floating horizon algorithm (WHAN, WOTI) can be used for the back face culling of the solid backbone model. Note that this process was not used as the ribbon model, since the the ribbon threads were simply avoided and destroyed regardless of their distance from a viewpoint and the viewers had to traverse the twist and width of the ribbon threads to observe the outer and width of the ribbon threads to observe a closer to observe the back of the backbone curve is closer to front.

One of the most important rationales behind our animation system is that it also leaves room for the display of the permum molecular models. As mentioned earlier, the spherical model itself is not animated in our system because protein engineers are more concerned about the conformational change represented by the backbone, which are usually hidden under Van der Waals surfaces in the spherical model representation.

However, the importance of the spherical model need not be overlooked since the geometry of spherical surfaces determines most of the crucial properties of the molecules. Often protein engineers need to map the interaction flow and query about the geometrical information of the atoms, to

a static view of a scene. In such a situation, the spherical model can be a useful tool for the analysis and thus is the most reason why we make the solid backbone model pass through the backbone atoms.

By making the backbone curve pass through the backbone atoms, viewers can locate the exact position of backbone atoms. Therefore, they may notify the anisotropic nature of their rate of contacts by clicking on the backbone atoms in the concerned area. The animation system, in turn, can display the spherical model of the backbone atoms in that specific area. Furthermore, extra atoms belonging to clicked residues can be displayed, provided with a proper data structure (Atoff, Ratoff) that could link the positional information of the atoms within a residue. Sincere sight rather plain nearby nitrogen and carbon in the same plane, specialized of a sequence of backbone atoms for right carbonyl can completely remove the structure of all atoms related to these backbone atoms with proper setup of a database (Atoff, Ratoff). This way, only those atoms which interest the viewer, which are hidden inside most of the time in the spherical representation of an entire molecule, can be revealed in detail.

In exchange for the pleasing visual appearance of approximating splines, which do not necessarily pass through the backbone atoms, the accuracy of the backbone curve can be enhanced. At least, backbone atoms lie on the backbone curve. Even if the backbone curve is made up of alpha carbons only, there is a greater chance that the nearby nitrogen and carbon atoms are to the backbone curve than in the case with the approximating splines. The approximating splines are best suited for the display of secondary structure patterns which are usually more in nature and thus tend to appear at the final stage of the folding process. Since the folding process in general is the transition from the random, unfolded conformation to the

secondary structural patterns, conventions states do not generally exhibit such patterned structures. Therefore the conformation at transition does not have to be visually pleasing. Indeed, the graphical representation should be focused on representing a state accurately the shape of the folding intermediates.

In practice, the property of passing through the backbone states acts as a bridge between the state and the dynamic display (i.e., animation). As long as there is a proper tool for baking the dynamic and static display, various state models presented previously can be explored and interpreted into the animation system. At the bottom the animation is interrupted for a static analysis, the system may go back to the nearest key frame and the static display can be started with the position of the alpha carbons on that key frame. Therefore, each static display can be labelled as part of the dynamic display with this scheme.

Figure 5.9 clarifies this point. In the upper left, a dot sphere model of a certain protein residue sequence is shown. The positions of the alpha carbon atoms are extracted by the joints of the straight line segments and the joints themselves are the centers of the dot spheres representing the alpha carbons. A smooth curve running along the dot spheres is an example of the approximating spline. Notice how the spline passes the centers of the alpha carbon atoms. The backbone curve for the situation at hand leaves the posture of the backbone atoms by a hinge margin. Therefore it is hard to track the position of backbone atoms on such a curve. The top right figure is a spherical model at the same molecule, and the bottom figure is a current view of a portion of that molecule. By reading the backbone curve pass through the backbone atoms, the position of the backbone atoms can be marked on the curve, and the alpha carbon locations can be computed by

referred to the animation system just by designating their positions $H_i$ on the other hand, the backbone curve is constrained by an approximating spline, there is no way to ensure the exact position of the backbone atoms upon mapping the secondary flow. Thus in the worse worse why the solid backbone model is set to pass through the backbone curves.

## Frame Interlocutor

The second aspect of the proposed system is the incorporation of inter-covering techniques into the solid backbone model. The interlocutory technique is a smooth area of computer graphics has been investigated in Chapter 4 in detail and thus this section will present its application to the context of the animation of protein folding.

To date, none of the existing molecular graphic systems has attempted to use the inter-covering technique in the context of animation. This is true even for the animation of the protein dynamics which can be loosely defined as conformational fluctuation over a given time span (Mait??, Mait?). Previous animations (Delhi, Kajer, Vanfelt) on molecular reformation was contain characteristics, as follows:

1. In the rendering aspect, most of the systems use the wireframe model as an animation tool

2. Since a full bond is represented in the wireframe model, high frequency motions of wireframe could distract the viewer's attention

3. Since it only the frame connecting the backbone atoms are displayed, the motion transition is adversely abrupt, since the joint angles of the wire frame can vary widely from one frame to another

By incorporating a smooth curve to the total backbone model, the problems associated with the first two characteristics of the telephone model could be removed. However, the last characteristic should be handled properly to obtain enhanced display of the protein folding process.

It is important to note that the steric influence of intermolecular backbone shapes is not sufficient for unwanted display, since the improper intramolecular backbone relies heavily on visual illusions. When a collection of beams representing the backbone is displayed sequentially, the visual illusion of saturation is possible only when a sufficient number of static frames, which are similar in shape, are provided. If the image of the current frame deviates remarkably from that of the previous frame, the previous image does not help saturate the scene, although it remains and becomes overlapped by the current frame visually.

In view of the preparation for the shapes of the folding intermediates, it can be said that visual recurrences between adjacent intermediates cannot be assured. Because of the limited technologies (SoftE), to date, only a few folding intermediates have been able to be captured and displayed. Theoretical postulations on the folding pathway (Chou4, Wood) does not necessarily document minute transitions from one frame to another.

Rather, a combination of the secondary structures, which folding in a very limited portion, is taken at the entire conformational change that folding process undergoes, is utilized well in the explanation of the folding pathway. For instance, in Figure 3-8, a formation of polypeptide chain is explained in term of the formation of a transient alpha helix.

Although a theoretical postulation can be described by such a schematic diagram, the animation of the folding transition involving a detailed conformational structure cannot be presented in such a discrete way. The

came holds true for calculation of the protein dynamics or energy minimization. For instance, if energy minimization is set to stop and produce the conformational data for a predefined small change in the energy level, the conformation thus produced can be varied greatly.

No matter what folding pathway is assumed, the sequent folding intermediates should be stored at key frames to exploit the visual illusion of animation using computer graphics technology.

Interweaving frames can prove to be a useful technique to compute an intermediary point by interpolating between two key drawings in animation. As was discussed in the Chapter 6, it smooths out the ultimate key frames by inserting extra frames between them so that it ensures smooth motion intuition.

By filling in a pair of chosen key frames with the inbetween frames, the resemblance between frames is exploited, so that the visual illusion leading to animation may take place. Given the conformations for a certain instant and the conformation at a later time, interweaving is the process of generating intermediary shapes. Figure 8-3 illustrates how the backbone curves change their shapes gradually by the incorporation of the inbetweening technique. Pairs of key frames (i.e., two folding intermediates of a protein structure) produced thus cannot ensure the smoothness of the motion. The discontinuity in the sequent folding intermediates cannot be solved in this process. For instance, there are two folding intermediates that cannot fully coalesce owing to kinematical limitation, resulting in gaps and discontinuities, as the conformation resulting from energy minimization. During playback of this sequence, the interpolating scheme computes these frames between key frames itself, regardless of whether there are one or two inbetween frames. This method contrasts with some of the previous approaches that specify a single frame number and compute it anymore, whereby the fixed contours of a single protein structure and converted it to more by a certain rule [Ras79].

Figure 9-7. An approximating spline used for the arbitration of the backbone curve.

unfolded
chain

Transient
α helices

Stabilised alpha
helices before
first folding unit

Figure 8.8 A schematic diagram showing the transition between folding states. Two segments of an unfolded polypeptide chain transiently become α-helical. These helices are then stabilised by the formation of a complex between the two segments.

key frame #1
(different frame #1)

key frame #2

inbetween
frame #1

inbetween
frame #2

inbetween
frame #3

inbetween
frame #4

inbetween
frame #5

Figure 5-1 A schematic diagram of a folding structure. Frame numbers are assigned in time sequence. Starting key frame (key frame #1) is treated as the first inbetween frame for reasons described in Chapter 4. A total of five inbetween frames are produced to smooth out the transition between the two key frames.

## System Design

This section describes a prototype system for folding estimation. The total backbone model described in the preceding section will be implemented and used as a basis learnt representing a folding intermediate. The substructuring techniques will be incorporated to produce extra tokens which would enable smooth transitions between two folding intermediates. In making up a solid backbone model, the G? dots force splines developed in Chapter 4 will be used to enable supplies-free representation of the backbone curve. The extra splits will be used in producing a wiggle-free motion path of an individual atom belonging to a folding intermediate over the given time span. The speed control technique developed in Chapter 4 will be used to generate a constant speed during the animation. Generated frames will be illustrated in the next section.

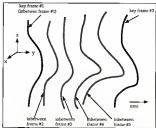The prototype system is made up of a combination of hardware: VAX 11/780, Evans/Sutherland PS/390 and a Macintosh Ilx with 8 megabytes of RAM and 128 megabytes of hard disk space. SVRT1, a molecular modeling software from Tripos Inc., was used for the preparation of the input data in our animation system.

Figure 9.18 shows a block diagram of the prototype system including the input preparation stage. The role of each block is as follows.

## 1. Input preparation

In the prototype system, we concentrate on the folding intermolds of a protein made up of twelve amino acid residues. Using Evans/Sutherland

79/390, a sequence of residues of a ubiquitin molecule which belongs to the alpha helical portion is extracted. It is:

K31-ILE32-GLN-VAL4-LEU-LYS-ALA40-

LYS37-ILE44-GLN45-LYS46-GLU48-GLY49,

where

ILE = isoleucine, GLU = glutamic acid, ASN = asparagine;

VAL = valine, LYS = lysine, ALA = alanine;

GLN = glutamine, ASP = aspartic acid, GLY = glycine.

The principal impetus behind the postorgic system, in terms of biochemical interest, is twofold. First, if we design and construct a protein by connecting sequences of each extracted residues in a geometrically linear fashion, can energy minimization yield the outer conformations (an alpha helical or one embedded in the original ubiquitin molecule). Second, is there any role of a specific one that characterizes the folding process at the designed molecule? The selection may suggest an answer to the stated question.

The VAX 11/780 was used to calculate the gradual, conformational change of the geometrically linear sequence of residues, upon energy minimization. In each run, the resulting conformational data with a certain energy level makes up a key frame to be used in the next steps. A total of fifty key frames were generated with successive reduction in the energy level. The output file format, MOL2 file format [SyVI6], is a collection of a three-dimensional Cartesian coordinates of all the atoms depicting the backbone structure of the protein.

## 2  Backbone contour

The reason why the backbone contour is expressed from the input preparation steps is to allow more flexibility of the system in that various

kernels of conformational data can be handled. Atomic coordinates from different sensors and different energy minimizations schemes can be incorporated into our system with a proper modification of its backbone sensor. The sensor outputs other *bit kernels* of alpha carbons or can be set to extract cubic backbone *items* such as *octagons*, *tetrons* and *alpha carbons* from the given conformational output file of the input stage.

## 3. Endpoint scenario

Generally, in the cubic sphere method, four data partitions are needed to evaluate a *mean* interpolating *test lower* data *power*. Two scenario *point* partitions are required to *initiate* and *maintain* the *sequential continuity* as was discussed *Chapter 5*. This is *true* both for the C7 and C7 *class* spline. In terms of the *data points* made up of the *twelve residue* regression, partitions of residues #3 and #4 are *used* to *provide* a *backbone curve* extrapolating between the residues #4 and #5. Therefore, *inside* the *sequential* data *points*, this *restriction* does *not cause* any *problems*. However, to *evaluate* an *interpolating* spline *running* from residues #1 and #2, an *extra data point* preceding residue #1 is *needed* and this data point *somehow alters* the curve *shape* in the *interpolated region*. The *same* holds *true* for the *evaluation* of the curve *running* from *residues* #11 and #12. Thus, *two extra data point* locations *should* be *supplied* as an *end medium*.

Although *some effort* has been *made* to *produce* a *possible end condition* (Boroff, Paris), an *additional complexity* of the *calculation* is *involved* in the *evaluation* of the *endpoint*. For *instance*, a *designed* *end condition* requires the *prescription* of *tangent vectors* at the *data points* and the *quadratic end condition* requires *evaluation* of the *second order derivatives*. However, no *generally-known best solution* to the *end condition*

problem exists until it is reasonable to assume a number of possible positions of the endpoints depending upon the structure encountered. In the prototype system, the endpoint positioning routine #1 is assumed to be in the apparent positions of residue #1 with respect to residue #1, and similarly for the endpoint following residue #13.

Figure 5-11 shows a data structure representing the positional information of the frames. The block composed of white squares represents the original data received from the input preparation steps. Each column of the block represents the position of eight carbons residing in a backbone conformation as input. For instance, the square labeled 1 in the three-dimensional Cartesian coordinate of the position of the alpha carbon atom belonging to the residue #1 carbons. Therefore, starting from the left, each column can be regarded as a sequence of key frames over the time interval and the total number of columns is 30, as was determined by the input preparation step.

While each column of Figure 5-11 shows a backbone conformation at a certain instant, a row in the figure represents the positional change of an alpha carbon over a given time interval. For instance, the first row of the block of white squares specifies the spatial history of the alpha carbon atom of the residue #1, and likewise the spatial history of the movement of the alpha carbons is encoded in the successive rows of the block.

The endpoint processor works in two ways in the figure; first, an additional structure residue position is attached to both ends of a column as shown in the figure, thereby making up each frame of 14 residues. These additional residues will be used as the endpoints for the generation of the backbone conformation in the spline interpolation module so that the interpolated curve may range from residue #1 to #13 in its entirety. Second, a pair of additional endpoints is generated for each row of the figure. These

endpoint are required to evaluate the motion path of a single atom in the time domain. Just as the endpoints are needed for drawing the cubic beziérline curve at a certain time, so are they needed for drawing the motion path of a single atom over a time interval. However additions all the extra atomic positions will generate two extra key frames, one preceding the first key frame, the other following the last key frame. The additional key frames will be used for the interpolations of the motion paths in the subsequent spline interpolation module.

## 6   Résolutions converter

This routine converts the collection of columns in Figure 5-31 into rows and vice versa (Since the first row of any spline interpolation is the interpolation of motion path in the time domain, and since the interpolation is applied for each atom, the column data should be converted into row data. Therefore, theorem this, each containing the locus of each residue, including this two extra endpoints, will be generated as a result of running the converter.

## 7   Spline interpolation module

The Incremental Knot Spacing method developed in Chapter 4, and the Q[2] four knot spline developed in Chapter 3 are applied to the motion speed interpolation and the backbone shape interpolation, respectively, in this module. Although two separate routines make up the module, they share a common routine in practice, since both adopt the Q[2] four knot spline as their interpolation scheme. The same interpolation method can be used both for the generation of the motion path of a single atom over a time interval, and

the generation of a backbone curve at a certain current. Consequently, the timing and the shape of inbetween frames are determined by this module.

The time reason why they are separately treated lies in the fact that the additional flexibility of controlling motion speed is required in the motion speed perception. Our prototype system is made to maintain constant speed between pairs of key frames, but one might need to slow down the change of the endomorture in a certain key frame range and this is possible by varying motion speed parameters in the module prob. beispekken module. By the same token, the number of points approximating a backbone curve can be arbitrarily chosen in the backbone steps interpolant. For instance, if a curve between two emotions is approximated by the combination of the line segments joining eight learned points, the curve can be said to be more accurately approximated than one constructed by joining four interval points. Since the number of emotion representing a large molecule could greatly be increased, there should be a trade-off between the number of learned data points and the computation time involved.

As a result of running the module prob interpolant, subsequent inbetween emote positions are generated for each row of Figure 9.11. For instance, if eight inbetween frames per pair of key frames are specified (or the motion path component, a total of 65 frames ($6 * 10 + 1$) are generated to render smooth transitions between backbone configurations. Since these interpolations can be made, creating up a row, and since the actual backbone conformation is built up upon the data points represented by columns, the output of the module prob interpolant should be fed into the row/column consumer.

Finally, the backbone shape interpolant is used to generate the backbone curve corresponding to each of the 401 frames which are produced by the motion path interpolant.

## 6 Transduce

There is a particular aspect of motion dynamics to be considered in connection with the initiation of protein folding. If the frames illustrated in Figure 5-12 are displayed sequentially, as in the global motion of the backbone curve will include the displacement of the backbone from one frame to another. The backbone curve, in its entirety, will be translated as if the backbone itself is moving. As a matter of fact, all the states involved in the energy minimization calculation are translated and the data given in the input preparation stage do contain the pertinent displacement.

Nevertheless, in the animation of the folding process, what we are interested in is the relative deformation of the inside conformation of the backbone curve and not the global translation of the curve. Suppose, for example, that the relative motion of the atoms and the earth is to be examined by an animation. If the scene were captured in such a way that the sun, and the circular movement of the moon and the earth around it, is included, viewers would be overwhelmed by the global motion of the moon and the earth around the sun. Although the relative motion between the moon and the earth is correctly activated, it might go unnoticed. By the same token, the relative inside deformation of the backbone curve may be ignored if the global displacement of the backbone curve is also included in the animation.

Therefore, a reduction point should be postulated on the backbone curve to facilitate the perception of the relative deformation of that curve

Once a reference point is determined on a certain frame, all the residue locations in the other frames should be translated such that they all meet at the reference point. Figure 5-12 shows the difference in the three animated frames resulting from variation of the reference point. The reference point on the first illustration is set to the N-terminus, while the second one is set to the C-terminus of an amino acid residue. If a reference point is fixed near the middle of the backbone curve, the animation will appear such that the backbone deformation propagates with respect to the reference point. In the prototype system, the position of the first residue is used as the reference point. However, it is important to note that no matter what strategy is taken, the relative position of an atom within a backbone curve remains fixed, since all of the atoms are translated as a group, by the translation module.

## 7 Display module

This module is responsible for the generation of the cylindrical surfaces running along the length of the backbone curve and there are some geometrical calculations involved in the animation of a cylinder. In practice, a cylinder can be approximated by an arbitrary number of polygons as shown in Figure 5-13. In this figure, four polygonal planes parallel to the line segment joining the center of surfaces A and B are drawn.

Suppose that the backbone curve shown in the figure is to be approximated by the line segments joining the base points $C_i$ through $C_j$. If we denote $a, b, c$ components of $C_i$ as

$$C_i = (x_i, y_i, z_i),$$

then the surface normal vector of the surface A can be approximated by

$$n = (a_2 - a_1) \wedge (a_3 - a_1), \quad |a_2 - a_1|$$

Therefore, a vector which lies on the surface and is also perpendicular to the surface normal can arbitrarily be written as,

$$v = (a_2 - a_1), \quad (a_2 - a_1), \quad 0)$$

Given the radius $R$ of the cylinder, a point $P$ on the node lying on the surface A can be calculated as,

$$
\begin{aligned}
P_x &= x_0 + R (q_0 - q_3) / v \mid \quad v \mid \\
P_y &= y_0 - R (q_0 - q_0) / \mid v \mid \\
P_z &= z_0
\end{aligned}
$$

where

$P_x, P_y, P_z$ are $x, y, z$ components of the point $P$.

Having set a point on the cavity, the rest of the points on that circle can be calculated easily by rotating the point with a certain angle with respect to the center of the cavity. However, in order to build a polygon, another point $P$ on the surface B is required to match the point $P$ on the surface A. A simple approach is taken in our prototype system. The matching point is assumed to be the intersection of the surface B, and the line parallel to the vector $(C_0 - C_0)$ while passing through the point $P$.

A perspective projection [Sut83, New79] was employed in the prototype system. If a square surrounding the backbone curve to the eye plane is assumed, the center of projection is made to be located at a certain distance from the center of the square with its direction parallel to the z-axis. Since the object further away from the center appears smaller in the perspective projection, the orthod provides a depth cue, an indication of which portions of the image correspond to parts of the object which are close or far away

A simple illumination model [Nabb, Regel0] is used to shade the solid landmass model such that the Lambertian cosine term is added to a constant diffuse light. According to the Lambert's cosine law, the intensity of a light reflected from a perfect diffusor is proportional to the cosine of the angle between the light direction and the surface normal. In addition, the ambient light explains the light scattered from the surroundings. Therefore, the light intensity of a certain polygon is

$$I = I_a \, k_a + I_l \, k_d \cos \theta$$

where

$I$ = the intensity of the light reflected from the polygon,

$I_a$ = the ambient radiance light intensity,

$k_a$ = the ambient diffuse reflection constant,

$I_l$ = the incident light intensity from the distant source,

$k_d$ = the diffuse reflection constant,

$\theta$ = the angle between the light direction and the surface normal

Since the light direction is assumed to be parallel in our system, and since the surface normal of a polygonal surface is constant, the reflected light intensity is constant for each polygon.

The display module is mainly concerned with the generation of state images of the individual frames. Various types of computionally-available sequencers can simply organize and display the individual frames for animation. However, we have not attented to this dissimulation to the generation of sequential frames based on the the solid landmass model approaches and further arrangement of the sequences is not covered.

Figure 5-12  Block diagram of the prototype system

Figure 9-11: Data structure of the animated frames used in the prototype system. Numbers represent the tendon sequence of the prototype system.

Figure 5-12. Referencing the frames in different places. If the M-terminus of the frame #1 is on to a reference point, all backbone atoms in the frame #2 and #3 are translated by the same amount as their M-termini are translated to the M-terminus of the frame #1.

surface B

a backbone curve

parallel distant light

$C_2$

surface A

surface normal

$C_1$

surface A.

Figure 5-63. A cylinder approximated by line polygons. The vector a represents the surface normal of the surface A. The vector i lies on 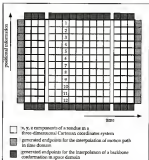the surface A and is perpendicular to the vector a. The interpolated backbone curve is approximated by a combination of the line segments joining the points $C_1$ from P and P are two matching points on surface A and B, respectively.

Some of the knives produced by the prototype system are illustrated in this section as visual experiments using the solid backbone model discussed above.

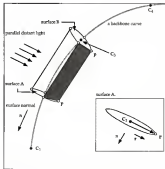Figure 5-4 shows the reshowing of the ubiquitous molecule shown in Figure 5-1. Instead of multiple threads, four polygonal surfaces are used to approximate the cylindrical surface representing the backbone in the solid backbone model. Although roughly approximating the surface and removal of the hidden surface is explicit enough to convey the spatial continuation of the backbone curve. The top figure shows the B-spline interpolation of the backbone curve by a combination of infills and the bottom figure shows the solid backbone model drawing. Although the solid backbone model is defined only in terms of the interpolating splines, this example is produced for the purpose of making a comparison between the two. As was shown, the B-spline is smoother than the line knife spline and is visually more pleasing. However, there appears to be more error in the B-spline approximation.

Figure 5-15 illustrates the alpha helical portion of the ubiquitin molecule shown in Figure 5-14. The helical portion ranges from residue 30 to residue 34 of the ubiquitin molecule. The top figure is the B-spline approximation and the bottom figure is the solid backbone model approximation. With white dots representing the location of the alpha carbon atoms, the B-splines can be said to be wide of the mark. Although aesthetically pleasing, the backbone curve is faulty far from the actual backbone atoms. In contrast, the solid backbone model drawing until achieves more accuracy by making the curve pass through all the positions of the alpha carbons. There is

a greater chance for the other backbone atoms, such as nitrogen and carbon, to reside in the vicinity of the carrer from it the smoothness is induced somewhat compared with that of the B-spline, usual continuity is preserved in this model.

Figure 5-4 shows a pair of nonomeric trates of the part of the alpha helix in Figure 5-42. The difference in the intensity of reflecting light is, shading can be shown clearly for individual polypeptid patches. Notice that the radius of the cylinder, and thus the sides of the polygon, is not regular because of the limited number of polygons used for the approximation. As the number of faes segments representing the backbone carva becomes, and as the number of polygons approximating a cylinder increases, each irregularity can be reduced.

Figure 5-37 is the array of the twelve residue sequences used as the input for energy minimation. It is modeled simply by concatenating each residue in a larer fashion. As a result of energy maximization 50 frames were available for our animation system. However, for purposes of the example, only the first three key frames are presented. Figure 5-18 illustrates a solid backbone model showing of the three key frames representing the folding structure of the protein. Notice that the conformational changes between adjacent frames are great enough to prevent visual libration when they are directly displayed for re animation.

Figures 5-39 and 5-40 illustrate the frames produced by the inbetweening technique. The motion speed is made constant by each joint belonging to the backbone cachctriene. Figure 5-39 shows the inbetween frames interpolating key frame #1 and #2, while Figure 5-40 shows three interpolating key frames #2 and #3. As can be seen, my two adjacent frames look quite similar so that the visual libration, leading to the animation can

happen if the frames are displayed sequentially in other words, accumulation of the small aerodynamic changes between adjacent frames makes the abrupt conformational changes of the key protons look smoother. For instance, although the first and last frame of Figure 9-18 are highly dissimilar, the individual conformational change in any pair of adjacent frames in the figure is very slight. This is the greatest advantage of the inbetweening technique.

## Summary

As an application of the free form splines developed in Chapter 3, and the speed control strategies developed in Chapter 4, a prototype animation system for protein folding has been presented. In addition, a new graphical model suitable for molecular animation has been proposed and tested in this chapter. The basic concepts involved in the solid backbone model are the solid drawing of the backbone and introduction of the inbetweening technique. In addition, the advantage of the interpolating spline employed in the solid backbone model over an approximating spline has been described and discussed.

Figure 9-16. Comparison of the Boφdale and fine form spheres. The backbone of the ubiquitin molecule discussed in the Figure 9-1 is shown as solid. Bottom figure is a representation of the solid backbone model and the top figure is a variation of the solid backbone model

Figure 9-13: An eighth helix portion of the ubiquitin molecule. The curve approximating the backbone of residues of strands in is drawn. White dots represent exact position of residue. Top figure is the B-spline approximation of the backbone, while the bottom one is the solid backbone model used in our prototype system.

Figure 5-16  Received view of a portion of the alpha helix in the figure 5-15 Top and bottom figures are zoomed in by a factor of 4 and 16, respectively.

Figure 5-17 Six linear arrangement of a sequence of twelve timelines used in the prototype system. Top left is the beginning Norarrensus and bottom right is the Cretaceous.

Figure 5-19 Three suspended backbone conformations used in the wavy domain in the prototype system. These frames and the coordinates are taken from the energy minimization output.

Figure 9-19. Six sets drawers produced by the interweaving technique. The second key frame is also shown for reference. The sequence is from top left to bottom left and from top right to bottom right. Notice the gradual change from one frame to another.

Figure 8-30. Six boxes produced by abbreviating key frames #2 and #5. The third key frame is also shown for reference. The sequence is from top left to bottom left and from top right to bottom right. Notice the gradual change from one boxes to another.

# CHAPTER 6
# CONCLUSIONS

This dissertation is primarily concerned with the anomalies of the porous folding process in biochemistry. The solid backbone model presented in this dissertation provides a convenient tool for the visualization of the more polypeptide chain in any animation environment.

More importantly, in the model, alpha carbons can be tabulated and marked as part of the backbone curve. When the animation abuses a stepped list further analysis of a tedial confirmation, each marked positions can be used as an interstice between system and viewer. Portions of a large macromolecule can be seen in the spherical model by designating the alpha carbons belonging to the area of interest, using a device such as a mouse. Provided with a proper database, the system can recognize the clicked location and seamlessly recover the location of other atoms belonging to the peptide planes centered at the alpha carbon.

As a result, those van der Waals surfaces buried under outside atomic surfaces can be exposed selectively. This was not possible in an animation relying entirely on the spherical model, despite the enormous amount of computation involved. In addition, by making the spline used in the solid backbone model pass directly through the alpha carbon locations, all of the static models developed earlier, such as the wireframe model and the solvent-accessible surface model, can be conveniently exploited as part of the animation system presented in this dissertation.

258

No matter what accounts for the determination of folding intermediates are used for the structures, the technique of interweaving is required to smooth out the motion sequence. Conformational data are continuously discrete whether the data are taken from hydrogen exchange, nuclear magnetic resonance, or energy minimization. The spatial pattern from one folding intermediate to another does not accurately exhibit a gradual transition. Therefore, the introduction of the interweaving technique is required to fill in gaps in the structure, since the visual chassis leading to the transition cannot be anticipated with such a discrete sequence of conformations. Despite its popularity in the context of computer animation, the interweaving technique has not been applied at any existing molecular animation system. This dissertation has developed a prototype molecular animation system incorporating the interweaving technique.

Splines are used in the prototype system in two ways: They are used for the generation of the curve passing through the backbone atoms and they are used for the generation of the motion path of individual atoms over a time span. In implementing these concepts, problems associated with conventional interpolating splines had to be resolved. Those interpolating splines with parametric continuity could not satisfy our needs because of the abnormal behavior called wiggles or loops. If the parametrically-continuous spline is employed, the backbone shape and the motion path cannot be free of wiggles, due to the unduly attained constraints applied to the spline. That is, at a certain data point, both the magnitude and direction of the left and right tangents should be identical. Another class of splines, known as visually continuous splines or GC2 however, encompasses a wider range of spline shapes while retaining the same smoothness visually as the parametrically-continuous spline.

Although an algorithmic approach to produce the curve was presented previously, no explicit formulation of the class of splines as yet exists. In previous approaches, recursive iteration based upon the construction of a Bézier polygon was mandatory to calculate the curve, and the complexity of this algorithm has limited popular use of this class of splines. The dissertation derives and proves an explicit matrix representation of the visually-continuous interpolating splines. The formulation was derived by allowing the insertion of variable knot parameters, with the same parabolic blending as was used for the formulation of the well known Overhauser curve. Based on mathematical representation, three characteristic types of the visually-continuous splines were developed.

An important aspect of the formulation of the splines presented in this dissertation is the recursive use of vector matrices. Exploiting the advantage of vectors, a simple graphical interface for the generation of a motion path was presented. An animator designing the motion path of an object can simply draw a line to denote the direction of the motion in corresponding data points. Such an extension of the visually-continuous splines may be applied to any application involving animation.

Another area of concern in the control of motion speed in the animation. A control strategy was required to ensure the proper motion speed of individual vectors when the sequence of folding intervolutions are displayed. Since the splines used in computer graphics are generally made to be parametric for certain purposes, the magnitude of the derivative on the curve, and thus the inter-frame distances, are parametric as well. Recent research in this area of animation suggested the possibility of controlling the speed using a means of calculation in the continuous domain involving estimation of the arc-gesture. However, in our prototype system, the folding

submacromolecules corresponding to key features must be part of the animation sequence. Our prototype system could not apply the continuous feature approach, since this technique does not necessarily list the key frames. A discrete domain alternative—the Incremental Knot Spacing method, combined with the Landmark Adjustment and Averaging Approximation—presented in this dissertation was shown to produce animations created of motion speed between frames. In addition, it is guaranteed to pass through the keyframe postures, so that the folding intermediaries can be treated and restored as part of the animation sequence.

The key contributions of this research, to computer graphics in general and molecular graphics in particular, can be summarized as follows:

1. Analytic mathematical formulation of visually-continuous splines, known as the $G^2$ class splines.

2. Development of a mathematical tool for the design of motion path in general animation systems, in an instance of the $G^2$ class splines.

3. Development of a method for the control of the motion speed in a spline-based animation system.

4. Development of a molecular model for the animated display of protein folding.

5. Introduction of the inbetweening technique in the field of molecular animation.

[Aho83]    Aho, A. V., Hopcroft, J. E., Ullman, J. D., *Data Structures and Algorithms*, 1983, Addison-Wesley, Reading, Massachusetts; New York, New York

[Ah079]    Akima, H., "A New Method of Interpolation and Smooth Curve Fitting Based on Local Procedures, *Journal of ACM*, vol 1, no 4, pp 580-602, 1979

[Ark83]    Arkwright, P. R., "A Sand-Site Holder: Surface Anatomical Procedures for Constructive Solid Geometry," *Computer Graphics*, vol 17, no 3, pp 59-67, 1983

[Bar80]    Baecker, R., "Interactive Computer-Mediated Animation," NAC-TR-61, Ph.D Thesis, Massachusetts Institute of Technology, 1969

[Bar84]    Barsky, B. A., "Exponential and Polynomial Methods for Applying Tension to an Interpolating Spline Curve," *Computer Vision, Graphics, and Image Processing*, vol 27, pp 1-18, 1984

[Bar84a]   Barsky, B. A., and Beatty, J. C., "Local Control of Bias and Tension in Beta-splines," *Computer Graphics*, vol 17, no 3, pp 193-218, 1983

[Bar84]    Barsky, B. A., and Bartel, P. L., *Computer Graphics and Geometric Modeling Using Beta-Splines*, 1988, Springer-Verlag, New York

[Bar84]    Bartels, R. H., Beatty, J. C., and Barsky, B. A., *An Introduction to Splines For Use in Computer Graphics and Geometric Modeling*, 1985, Morgan Kaufmann Publisher, Los Altos, California

[Bar81]    Barsky, R. H., and Hamilton, J., "Speed Adjustment for Key-Frame Interpolation," *Graphics Interface 91*, pp 14-19, 1981

[Bar79]    Bézier, P. E., *Numerical Control Mathematics and Applications*, 1972, John Wiley & Sons, New York

263

[Bal76] Balzer, P., "Mathematical and Practical Possibilities of INFERS," Computer Aided Geometric Design, Barnhill and Riesenfeld (eds.), 1974, Academic Press, New York.

[Bou77] Bouwhuis, J. E., "An Incremental Algorithm for Digital Display of Circular Arcs," Communication of ACM, vol 20, n.4, pp 238-246, 1977.

[Bow77] Bowyer, A. and Anderssen, D. C., "Visual Interaction With Overhanser Curves and Surfaces," Computer Graphics, vol 14, n.3, pp 131-139, 1977.

[Bro85] Brooks, F. B., Bourghow, R. E., Obonere, B. D., Omera, D. J., Swentdonohoe, S., and Kopfaw, M., "CHARISMA: A Program for Macromolecular Energy Minimization, and Dynamics Calculations," Journal of Computational Chemistry, vol 6, n.2, pp 187-217, 1985

[Ber77] Bartnyk, N., and Weir, M., "Computer-Generated Key-frame Animation," Journal of SMPTE, vol.80, pp 149-153, 1971

[Ber71] Bartnyk, N., and Weir, M., "Interactive Skeleton Techniques for Enhancing Motion Dynamics in Key Frame Animation," Communication of ACM, vol.19, n.10, pp 564-569, 1976

[Ca76] Caxton, M. P., Differential Geometry of Curves and Surfaces, 1976, Prentice-Hall, Englewood Cliffs, New Jersey.

[Ca79] Casson, M., and Begg, C. E., "Algorithm for Motion Models of Proteins," Journal of Molecular Graphics, vol.4, no.4, pp 157-162, 1986.

[Co84] Cantrell, E., and Rose, B., "A Class of Local Interpolating Splines," Computer Aided Geometric Design, Barnhill and Riesenfeld (eds.), 1974, Academic Press, New York.

[Ca78] Cantrell, E., "The Problems of Computer-Assisted Animation," Computer Graphics, vol 12, no.3, pp 348-353, 1978.

[Ca69] Cheville, J., Vernon, W. C., Monnet, T., Brenner, B., Gelin, and Honrana, G. M., "MANOSOC: A Graphics Program for Analyzing and Modeling Molecular Structure and Proximes," Journal of Molecular Graphics, vol 4, pp 159-162, 1988

Clto84   Clothiau, C., "Principles that Determine the Structure of Proteins," Annual Review of Biochemistry, vol. 53, pp 537-572, 1984.

Clto87   Clout. F. M., Miknt. P. E. and Jou, W., "Quantitative Reappraisal of Covered Expressions for Molecular Forms Binding in Substrate-Induce Chromatography," Biophysical Chemistry, vol.50, pp 151-161, 1974.

Cls85    Cluss, P. W. and Jou, W., "Molecular Conformation of Ubiquitinated Structures and the Implications for Regulating Structures," Journal of Molecular Biology, 1985.

Cls86    Clerk, M., Courae H. E. O. and Oydenkrach, N. H., "Validation of the General Purpose Tripos 5.2 Force Field," Journal of Computational Chemistry, vol.13, no.8, pp 982-1012, 1989

Coo89    Connolly, M. L., "Solvent-Accessible Surfaces of Proteins and Nucleic Acids," Science, vol.221, no.601, pp 709-713, 1983.

Coo74    Coorn, S. A., "Section Patties and Rolylatt Curves," Computer Aided Geometric Design, Barnhill and Riesenfeld (eds.), 1974, Academic Press, New York.

Crw79    Crow, F. C., "The Aliasing Problem in Computer Generated Shaded Images," Communication of ACM, vol.20, no.11, pp 799-805, 1977

Cus83    Crow, F. C., "A Comparison of Antialiasing Techniques," IEEE Computer Graphics and Applications, vol.1, no.1, pp 40-50, 1981.

Daa79    Davu, C. J., An Introduction to Solitary Surface Volume 1, 1984, Addison-Wesley, New York.

deH79    deHoor, C., A Pastoral Guide to Splines, 1978, Springer-Verlag, New York.

Dee82    DeNobe, P., Fantbun M., and Wodak, S., "Interactive Computer Animation of Macromolecules," Journal of Molecular Graphics, vol.1, no.3, pp 103-106, 1984.

Doo79    Doree, T. D., and Bartley, F. A., "Geometric Continuity, Shape Parameters, and Geometric Constructions for Catmull-Rom Splines," ACM Transactions on Graphics, vol.7, no.1, pp 1-41, 1988.

[Elm89]  Elmasri, R., and Navathe, S. B., *Fundamentals of Database Systems*, 1989, Benjamin/Cummings Publishing Company, California.

[Far88]  Farin, G. "Some Remarks on $v^2$ Splines", Computer Aided Geometric Design, vol.5, no3, pp 309-316, 1988.

[Far90]  Farin, G., *Curves and Surfaces for Computer Aided Geometric Design*, 1990, Academic Press, New York.

[Fed83]  Feldmann, R. J., and Levitt, M., "Molecular Dynamics of Bovine Pancreatic Trypsin Inhibitor," 1983, Film distributed by Association of American Medical Colleges as Study of the Indices Cue at Bovine Trypsin,", Proceedings of National Academy of Science, vol.76, no.11, pp.5828-5632, 1979.

[Fod84]  Fortenick, R., and Zeller, M., *Current Communications in Molecular Biology*, 1984, Cold Spring Harbor Laboratory, New York.

[Fol90]  Foley, J. D., and Van Dam, A., *Fundamentals of Interactive Computer Graphics*, 1990, Addison-Wesley, Reading, Massachusetts, New York.

[Fol84]  Foley, J. A., "Local Control of Interval Tension Using Weighted Splines," Computer Aided Geometric Design, vol.3, pp.281- 294.

[For77]  Forrest, A. R., "Interactive Interpolation and Approximation by Bezier Polynomials," *Computer Journal*, vol.15, no.1, pp.71-79, 1972.

[Fra81]  Franklin, W. R., "On Convex and Other Methods For The Representation of Curved Surfaces," Computer Graphics and Image Processing, vol.1, pp 261-296, 1972.

[Fra88]  Franklin, W. R., "An Exact Hidden Sphere Algorithm That Operates in Linear Time," Computer Graphics and Image Processing, vol.18, pp.366-379, 1980.

[Fu87]  Fu, K. S., Gonzalez, R. C., and Lee, C. S. *Robotics: Control, Sensing, Vision, and Intelligence*, 1987, McGraw-Hill, New York.

Fe84       Fejan, T., Benoson, S. and Meyer Z. "Icosahedon-ar Representations in the Active Site of Porcine Pancreatic Kemoxe with Acetyl-alanine-prolinechymotrypsin by Means of Molecular Dynamics Calculation," Journal of Molecular Graphics, vol 4, no 4, pp 208-212, 1986.

Ga85       Gisha, C, Sun, I., Protea Folding, Academic Press, New York.

Ga81       Glassner, A. S., "An Overview of Ray Tracing," Ray Tracing, Glassner, A. S. (ed.), 1989, Academic Press, San Diego, California.

Gae81      Gordon, W. J., and Rosenfeld, R. F., "Bernstein-Bezier Methods for the Computer-Aided Design of Free-form Curves and Surfaces," Journal of ACM, vol 21, no 2, pp 293-210, 1974.

Gae81a     Gordon, W. J., and Rosenfeld, R. F., "B-Spline Curves and Surfaces," Computer-aided Geometric Design, Barnhill and Riesenfeld (ed.), 1974, Academic Press, New York.

Ha67       Harrington, S., Computer Graphics, 1983, McGraw-Hill, New York.

He86       Heam, D., and Baker, M. P., Computer Graphics, 1986, Prentice-Hall, Englewood Cliffs, New Jersey.

He84       Herbison-Evans, D., "Nudy 2: A Humanic Utility Employing Ellipsoid Solids," Personal, Computer Graphics, vol 18, no 3, pp 354-356, 1978.

Ja74       Jacques, J., A First Course in Computing and Numerical Methods, 1970, Addison-Wesley, Reading, Massachusetts, New York.

Jae81      Jacques, J., and Smit, T., Numerical Analysis, 1997, Chapman and Hall, New York.

Ja75       Jan, W., and Chon, P. W., "Molecular Mechanics of the Extraction of Chobe Acid Micelles," Journal of Molecular Graphics, in press.

Ki82       Kim, P. S., and Baldwin, R. L., "Specific Intermediates in the Folding Reactions of Small Proteins and the Mechanism of Protein Folding," Annual Review of Biochemistry, vol 51, pp 459-489, 1982.

Kno67    Knowlton, K., "Computer-Aided Definition, Manipulation, and Depiction of Objects Composed of Spheres," *Computer Graphics*, vol.15, no.1, pp.10-71, 1981.

Kro86    Krichevets, D.H. U., and Koride, A.H., "Interpolating Splines with Local Tension, Continuity and Bias Control," *Computer Graphics*, vol.18, no.3, pp.33-41, 1984.

Lan82    Lane, J. M., Carpenter, L. C., Whitted, J. T., and Blinn, J. F., "Scan Line Methods Displaying Parametrically Defined Surfaces," *Communications of ACM*, vol.23, no.1, pp.23-24, 1980.

Lee85    Lancaster, P., and Šalkauskas, K., *Curve and Surface Fitting*, Academic Press, London.

Las87    Lasseter, J., "Principles of Traditional Animation Applied to 3D Computer Animation," *Computer Graphics*, vol.10, no.4, pp.20-34, 1987.

Lev84    Levinthal, C., "Molecular Model-building by Computer," *Scientific American*, vol.214, no.6, pp.42-52, 1966.

Lev88    Levitt, M., "Protein Folding by Restrained Energy Minimization and Molecular Dynamics," *Journal of Molecular Biology*, vol.170, pp.723-764, 1983.

Lev83    Levitt, M., "Molecular Dynamics of Native Protein I. Computer Simulation of Trajectories," *Journal of Molecular Biology*, vol.168, pp.595-620, 1983.

Lev84    Levitt, M., and Sharon, R., "Differential Geometry of Proteins," *Journal of Molecular Biology*, vol.181, pp.145-162, 1985.

Man78    Manning, J. R., "Continuity Conditions for Spline Curves," *The Computer Journal*, vol.17, no.2, pp.181-186, 1974.

Man80    Montka, C. M., "Parameterization by Arc Length and Geometrical Considerations," *Computer-Aided Design*, vol.10, no.1, pp.23-26, 1980.

Max82    Max, N., "A COBALL-screen with Shading and Highlights," *Computer Graphics*, vol.13, no.2, pp.165-173, 1979.

Max82    Max, N. L., "Computer Representation of Molecular Surfaces," *IEEE Computer Graphics and Applications*, vol.3, pp.21-29, August, 1983.

[Mon88] Moo, N. J., and Gossol, E. D., "Spheroidal Parametric Molecular Surfaces," *IEEE Computer Graphics and Applications*, vol 8, no 4, pp 52-50, July, 1988.

[Mor77] Morales, D., *Fachematics, 277*, Academic Press, New York

[MR81] Miller, D. E., Jea, M., Horemann, B. P., and Choo, F. M., "The Effect of Kill-Sail on Cochrane Antmyloma," *Unpublished, in press.*

[Nar79] Newman, W. M., Sproull, R. F., *Principles of Interactive Computer Graphics*, 1979, McGraw Hill, New York.

[Net74] Niebson, G. M., "Some Provaliser Polynomial Alternatives to Splines Under Tension," *Computer Aided Geometric Design*, Barnhill and Boreswhld (eds.), 1974, Academic Press, New York.

[Oso68] Oreshauser, A. W., *Analysis Definition of Curves and Surfaces by Parabolic Blasting, 1968*, Scientific Research Dtvd Publication, Ford Motor Company

[Pul87] Pacion, T. C., and Handman, F. H., "Corners Over Spheres, A New Method for Rapid CPK Image Generation," *Journal of Molecular Graphics*, vol 6, pp 149-156, 1988.

[Pqr78] Payer, M. E., "Technical Trends in Molecular Graphics," in *Current Communications in Molecular Biology*, 1984, Cold Spring Harbor Laboratory, New York.

[Pav81] Partez, F. R., "Spheroid Shading," *Computer Graphics*, vol 12, no 3, pp 282-285, 1978.

[Pow81] Press, W. H., Harvey, B. P., Teukoffsky, S. B., and Veterling, W. T., *Numerical Recipes in C, 1986, Cambridge University Press, New York*

[Rav81] Rawes, W., "Referencing for Computer Animation Uttlizing Moving-Point Constraints," *Computer Graphics*, vol 15, no 2, pp 263-270, 1981.

[Ric81] Richardson, J. S., "The Anatomy and Taxonomy of Protein Structure," *Advanced Protein Chemistry*, vol 34, pp 167-339, 1981.

[Ric81] Rocheschi, F. M., "The Protein Folding Problem," *Scientific American*, pp 54-63, Jan, 1983.

[Rog76] Rogers, D. F. and Adams, J. A., *Mathematical Elements for Computer Graphics*, 1976, McGraw-Hill, New York.

[Ros66] Rosenman, J. G., and Argyrs, P., "Proton Folding," *Annual Review of Biochemistry*, vol 66, pp 497-502, 1981.

[Sch83] Salkauskas, K., "C¹ Splines for Interpolation of Rapidly Varying Data," *Rocky Mountain Journal of Mathematics*, vol 14, pp 239-250, 1984.

[Sch84] Sarraf, H., *The Design and Analysis of Spatial Data Structures*, 1990, Addison-Wesley, Reading, Massachusetts, New York.

[Sch86] Schweikert, D. G., "An Interpolation Curve Using a Spline in Tension," *Journal of Math. Phys.*, vol 45, pp 312-317, 1966.

[Sch86a] Schulz, R. H., Inkinen, R. H., *Principles of Protein Structure*, 1979, Springer-Verlag, New York.

[Sm83] Smith, A. R., "Spline Tutorial Notes," *ACM SIGGRAPH '84 Course Notes: Introduction to Computer Animation*, pp 131-172, 1983.

[Sm86] Smulicowner, I., "On Display of Space Filling Atomic Models in Real-Time," *Computer Graphics*, vol 12, no 3, pp 157-172, 1978.

[Su89] Smulicowner, I., Khumaru, A. S., "Display of Molecular Models with Interactive Graphics," *IEEE Computer Graphics and Applications*, vol 4, no 1, pp 24-34, 1988.

[Su88] Smulicowner, J., Lrinoke, P. V., Wong, S. L., and Zhou, X., "Molecular Structure Visualization and Optimal Surface Reconstruction," *NCGA Conference Proceedings*, pp 169-177, March, 1989.

[Sz87] Szeliski, P. H., and Padfor, H., "Parametric Surface Interpolation Incorporating Kinetic Adjustment and Blending Curves," *Computer Graphics*, vol 16, no 3, pp 255-263, 1985.

[Sty42] Szyrej, L., Farskievitz, 1983, Willi Freeman And Company, New York.

[Te88] Thalmann, M. M., and Thalmann, D., *Computer Animation: Theory and Practice*, 1985, Springer-Verlag, Tokyo.

Tha85    Thalmann, N. M., and Thalmann, D. *New Trends in Computer
         Graphics: Proceedings of Computer Graphics International '88*,
         1988, Springer-Verlag, New York.

Tha88    Thomas, F., and Johnstone, O. *Disney Animation: The
         Illusion of Life*, 1981, Abbeville Press, New York.

Tod83    Todd, S., and Collett, J. "Animation in the Winchester
         Graphics System," *Journal of Molecular Graphics*, vol. 1, no. 2,
         pp.56-60, 1983.

Tri81    "TRYPL," Molecular Modeling Software, Harvey Mudd, Tripos
         Associates, a Subsidiary of Evans and Sutherland, St. Louis,
         Missouri, 1985.

Voo89    Voss, D. and Kurt, J. G. *Bookkeeping*, 1990, John Wiley & Sons,
         New York.

Wha84    Wang, E. C. "On Visibility Determination in Surfaces with
         Periodic Colored Dots," *Ph.D Thesis, Univ. of Florida*, 1980

Whi80    Whitted, J. T. "An Improved Illumination Model for Shaded
         Display," *Communication of ACM*, vol.23, no.6, pp.343-349, 1980

Wil80    Wright, F. S., Dineen, H. J. and Lerner, R. A. "Co-formation of
         Peptide Fragments of Proteins as Approach Relevance
         Implications for Initiation of Protein Folding," *Biochemistry*,
         vol.27, no.30, pp.7167-7175, 1988

Wri83    Wright, T. J. "A Two-Space Solution to the Hidden Line
         Problem for Plotting Functions of Two Variables," *IEEE
         Transactions on Computers*, vol.C-22, pp.28-33, 1973.

# BIOGRAPHICAL SKETCH

Moonsok Jou was born in Pohang, Korea, on November 7, 1959. He received his B.S. degree in electronic engineering from Seoul National University, Seoul, Korea, in 1983. After graduation, he was employed as a system engineer by International Business Machines, Incorporated, Seoul, Korea. His duties included instruction on IMS database, SNA concepts, and Systat/370 Assembler Language.

He was admitted to the master's program in electrical engineering at the University of Florida in the fall semester of 1985 and worked as a research assistant at the Database Research Center, Department of Computer and Information Science. He received his M.S. degree in electrical engineering in the fall semester of 1987. With the advice of Dr. John Staudhammer and Dr. Paul W. Chun, he continued on a Ph.D. program at the University of Florida. He pursed the doctoral qualifying examination in the fall semester, 1989.

His current research interests are computer graphics image processing, database design, digital signal processing, and VLSI design areas.

I certify that I have read this study and that in my opinion it conforms to acceptable standards of scholarly presentation and is fully adequate, in scope and quality, as a dissertation for the degree of Doctor of Philosophy.

John Straubhaamer, Chairman
Professor of Electrical Engineering

I certify that I have read this study and that in my opinion it conforms to acceptable standards of scholarly presentation and is fully adequate, in scope and quality, as a dissertation for the degree of Doctor of Philosophy.

Pang W. Chen
Professor of Biochemistry and Molecular
Biology

I certify that I have read this study and that in my opinion it conforms to acceptable standards of scholarly presentation and is fully adequate, in scope and quality, as a dissertation for the degree of Doctor of Philosophy.

Jack K. Solis
Professor of Electrical Engineering

I certify that I have read this study and that in my opinion it conforms to acceptable standards of scholarly presentation and is fully adequate, in scope and quality, as a dissertation for the degree of Doctor of Philosophy.

Isaac Norton
Assistant Professor of Electrical Engineering

I certify that I have read this study and that in my opinion it conforms to acceptable standards of scholarly presentation and is fully adequate, in scope and quality, as a dissertation for the degree of Doctor of Philosophy.

Chang W. Park
Associate Professor of Chemical Engineering

This dissertation was submitted to the Graduate Faculty of the College of Engineering and to the Graduate School and was accepted as partial fulfillment of the requirements for the degree of Doctor of Philosophy.

August, 1995

Winfred M. Phillips
Dean, College of Engineering

Madelyn M. Lockhart
Dean, Graduate School